
Sample-Efficient Reinforcement Learning for Linearly-Parameterized MDPs with a Generative Model

Bingyan Wang*
Princeton University
bingyanw@princeton.edu

Yuling Yan*
Princeton University
yulingy@princeton.edu

Jianqing Fan
Princeton University
jqfan@princeton.edu

Abstract

The curse of dimensionality is a widely known issue in reinforcement learning (RL). In the tabular setting where the state space \mathcal{S} and the action space \mathcal{A} are both finite, to obtain a nearly optimal policy with sampling access to a generative model, the minimax-optimal sample complexity scales linearly with $|\mathcal{S}| \times |\mathcal{A}|$, which can be prohibitively large when \mathcal{S} or \mathcal{A} is large. This paper considers a Markov decision process (MDP) that admits a set of state-action features, which can linearly express (or approximate) its probability transition kernel. We show that a model-based approach (resp. Q-learning) provably learns an ε -optimal policy (resp. Q-function) with high probability as soon as the sample size exceeds the order of $\frac{K}{(1-\gamma)^3 \varepsilon^2}$ (resp. $\frac{K}{(1-\gamma)^4 \varepsilon^2}$), up to some logarithmic factor. Here K is the feature dimension and $\gamma \in (0, 1)$ is the discount factor of the MDP. The results is applicable to the tabular MDPs by taking the coordinate basis with $K = |\mathcal{S}| \times |\mathcal{A}|$. Both sample complexity bounds are provably tight, and our result for the model-based approach matches the minimax lower bound. Our results show that for arbitrarily large-scale MDP, both the model-based approach and Q-learning are sample-efficient when K is relatively small, and hence the title of this paper.

1 Introduction

Reinforcement learning (RL) studies the problem of learning and decision making in a Markov decision process (MDP). Recent years have seen exciting progress in applications of RL in real world decision-making problems such as AlphaGo [46, 47] and autonomous driving [33]. Specifically, the goal of RL is to search for an optimal policy that maximizes the cumulative reward, based on sequential noisy data. There are two popular approaches to RL: model-based and model-free ones.

- The model-based approaches start with formulating an empirical MDP by learning the probability transition model from the collected data samples, and then estimating the optimal policy / value function based on the empirical MDP.
- The model-free approaches (e.g. Q-learning) learn the optimal policy or the optimal (action-)value function from samples. As its name suggests, model-free approaches do not attempt to learn the model explicitly.

*Equal contribution.

Generally speaking, model-based approaches enjoy great flexibility since after the transition model is learned in the first place, it can then be applied to any other problems without touching the raw data samples. In comparison, model-free methods, due to its online nature, are usually memory-efficient and can interact with the environment and update the estimate on the fly.

This paper is devoted to investigating the sample efficiency of both model-based RL and Q-learning (arguably one of the most commonly adopted model-free RL algorithms). It is well known that MDPs suffer from the curse of dimensionality. For example, in the tabular setting where the state space \mathcal{S} and the action space \mathcal{A} are both finite, to obtain a near optimal policy or value function given sampling access to a generative model, the minimax optimal sample complexity scales linearly with $|\mathcal{S}| \times |\mathcal{A}|$ [5, 2]. However contemporary applications of RL often encounters environments with exceedingly large state and action spaces, whilst the data collection might be expensive or even high-stake. This suggests a large gap between the theoretical findings and practical decision-making problems where $|\mathcal{S}|$ and $|\mathcal{A}|$ are large or even infinite.

To close the aforementioned theory-practice gap, one natural idea is to impose certain structural assumption on the MDP. In this paper we follow the feature-based linear transition model studied in [63], where each state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$ admits a K dimensional feature vector $\phi(s, a) \in \mathbb{R}^K$ that expresses the transition dynamics $\mathbb{P}(\cdot|s, a) = \Psi\phi(s, a)$ for some unknown matrix $\Psi \in \mathbb{R}^{|\mathcal{S}| \times K}$ which is common for all (s, a) . This model encompasses both the tabular case and the homogeneous model in which the state space can be partitioned into K equivalent classes. Assuming access to a generative model [31, 32], under this structural assumption, this paper aims to answer the following two questions:

How many samples are needed for model-based RL and Q-learning to learn an optimal policy under the feature-based linear transition model?

In what follows, we will show that the answer to this question scales linearly with the dimension of the feature space K and is independent of $|\mathcal{S}|$ and $|\mathcal{A}|$ under the feature-based linear transition model. With the aid of this structural assumption, model-based RL and Q-learning becomes significantly more sample-efficient than that in the tabular setting.

Our contributions. We focus our attention on an infinite horizon MDP with discount factor $\gamma \in (0, 1)$. We use ε -optimal policy to indicate the policy whose expected discounted cumulative rewards are ε close to the optimal value of the MDP. Our contributions are two-fold:

- We demonstrate that model-based RL provably learns an ε -optimal policy by performing planning based on an empirical MDP constructed from a total number of

$$\tilde{O}\left(\frac{K}{(1-\gamma)^3 \varepsilon^2}\right)$$

samples, for all $\varepsilon \in (0, (1-\gamma)^{-1/2}]$. Here $\tilde{O}(\cdot)$ hides logarithmic factors compared to the usual $O(\cdot)$ notation. To the best of our knowledge, this is the first theoretical guarantee for model-based RL under the feature-based linear transition model. This sample complexity bound matches the minimax limit established in [63] up to logarithmic factor.

- We also show that Q-learning provably finds an entrywise ε -optimal Q-function using a total number of

$$\tilde{O}\left(\frac{K}{(1-\gamma)^4 \varepsilon^2}\right)$$

samples, for all $\varepsilon \in (0, 1]$. This sample complexity upper bound improves the state-of-the-art result in [63] and the dependency on the effective horizon $(1-\gamma)^{-4}$ is sharp in view of [34].

These results taken collectively show the minimax optimality of model-based RL and the sub-optimality of Q-learning in sample complexity.

2 Problem formulation

This paper focuses on tabular MDPs in the discounted infinite-horizon setting [10]. Here and throughout, $\Delta_{d-1} := \{v \in \mathbb{R}^d : \sum_{i=1}^d v_i = 1, v_i \geq 0, \forall i \in [d]\}$ stands for the d -dimensional probability simplex and $[N] := \{1, 2, \dots, N\}$ for any $N \in \mathbb{N}^+$.

Discounted infinite-horizon MDPs. Denote a discounted infinite-horizon MDP by a tuple $M = (\mathcal{S}, \mathcal{A}, P, r, \gamma)$, where $\mathcal{S} = \{1, \dots, |\mathcal{S}|\}$ is a finite set of states, $\mathcal{A} = \{1, \dots, |\mathcal{A}|\}$ is a finite set of actions, $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{|\mathcal{S}|-1}$ represents the probability transition kernel where $P(s'|s, a)$ denotes the probability of transiting from state s to state s' when action a is taken, $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ denotes the reward function where $r(s, a)$ is the instantaneous reward received when taking action $a \in \mathcal{A}$ while in state $s \in \mathcal{S}$, and $\gamma \in (0, 1)$ is the discount factor.

Value function and Q-function. Recall that the goal of RL is to learn a policy that maximizes the cumulative reward, which corresponds to value functions or Q-functions in the corresponding MDP. For a deterministic policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ and a starting state $s \in \mathcal{S}$, we define the value function as

$$V^\pi(s) := \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k r(s_k, a_k) \mid s_0 = s \right]$$

for all $s \in \mathcal{S}$. Here, the trajectory is generated by $a_k = \pi(s_k)$ and $s_{k+1} \sim P(s_{k+1}|s_k, a_k)$ for every $k \geq 0$. This function measures the expected discounted cumulative reward received on the trajectory $\{(s_k, a_k)\}_{k \geq 0}$ and the expectation is taken with respect to the randomness of the transitions $s_{k+1} \sim P(\cdot|s_k, a_k)$ on the trajectory. Recall that the immediate rewards lie in $[0, 1]$, it is easy to derive that $0 \leq V^\pi(s) \leq \frac{1}{1-\gamma}$ for any policy π and state s . Accordingly, we define the Q-function for policy π as

$$Q^\pi(s, a) := \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k r(s_k, a_k) \mid s_0 = s, a_0 = a \right]$$

for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. Here, the actions are chosen by the policy π except for the initial state (i.e. $a_k = \pi(s_k)$ for all $k \geq 1$). Similar to the value function, we can easily check that $0 \leq Q^\pi(s, a) \leq \frac{1}{1-\gamma}$ for any π and (s, a) . To maximize the value function or Q function, previous literature [9, 51] establishes that there exists an optimal policy π^* which simultaneously maximizes $V^\pi(s)$ (resp. $Q^\pi(s, a)$) for all $s \in \mathcal{S}$ (resp. $(s, a) \in \mathcal{S} \times \mathcal{A}$). We define the optimal value function V^* and optimal Q-function Q^* respectively as

$$V^*(s) := \max_{\pi} V^\pi(s) = V^{\pi^*}(s), \quad Q^*(s, a) := \max_{\pi} Q^\pi(s, a) = Q^{\pi^*}(s, a)$$

for any state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$.

Linear transition model. Given a set of K feature functions $\phi_1, \phi_2, \dots, \phi_K : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, we define ϕ to be a feature mapping from $\mathcal{S} \times \mathcal{A}$ to \mathbb{R}^K such that

$$\phi(s, a) = [\phi_1(s, a), \dots, \phi_K(s, a)] \in \mathbb{R}^K.$$

Then we are ready to define the linear transition model [63] as follows.

Definition 1 (Linear transition model). *Given a discounted infinite-horizon MDP $M = (\mathcal{S}, \mathcal{A}, P, r, \gamma)$ and a feature mapping $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^K$, M admits the linear transition model if there exists some (unknown) functions $\psi_1, \dots, \psi_K : \mathcal{S} \rightarrow \mathbb{R}$, such that*

$$P(s'|s, a) = \sum_{k=1}^K \phi_k(s, a) \psi_k(s') \quad (1)$$

for every $(s, a) \in \mathcal{S} \times \mathcal{A}$ and $s' \in \mathcal{S}$.

Readers familiar with linear MDP literatures might immediately recognize that the above definition is the same as the structure imposed on the probability transition kernel P in the linear MDP model [63, 30, 65, 26, 52, 56, 60]. However unlike linear MDP which also requires the reward function $r(s, a)$ to be linear in the feature mapping $\phi(s, a)$, here we do not impose any structural assumption on the reward.

Example 1 (Tabular MDP). *Each tabular MDP can be viewed as a linear transition model with feature mapping $\phi(s, a) = e_{(s,a)} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$ (i.e. the vector with all entries equal to 0 but the one corresponding to (s, a) equals to 1) for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. To see this, we can check that Definition 1 is satisfied with $K = |\mathcal{S}| \times |\mathcal{A}|$ and $\psi_{(s,a)}(s') = \mathbb{P}(s'|s, a)$ for each $s, s' \in \mathcal{S}$ and $a \in \mathcal{A}$. This example is a sanity check of Definition 1, which also shows that our results (Theorem 1 and 2) can recover previous results on tabular MDP [2, 34] by taking $K = |\mathcal{S}| \times |\mathcal{A}|$.*

Example 2 (Simplex Feature Space). *If all feature vectors $\{\phi(s, a)\}_{(s,a) \in \mathcal{S} \times \mathcal{A}}$ fall in the probability simplex Δ_{K-1} , a linear transition model can be constructed by taking $\psi_k(\cdot)$ to be any probability measure over \mathcal{S} for all $k \in [K]$.*

A key observation is that the model size of linear transition model with known feature mapping ϕ is $|\mathcal{S}|K$ (the number of coefficients $\psi_k(s')$ in (1)), which is still large when the state space \mathcal{S} is large. In contrast, it will be established later that to learn a near-optimal policy or Q-function, we only need a much smaller number of samples, which depends linearly on K and is independent of $|\mathcal{S}|$.

Next, we introduce a critical assumption employed in prior literature [63, 67, 45].

Assumption 1 (Anchor state-action pairs). *Assume there exists a set of anchor state-action pairs $\mathcal{K} \subset \mathcal{S} \times \mathcal{A}$ with $|\mathcal{K}| = K^2$ such that for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, its corresponding feature vector can be expressed as a convex combination of the feature vectors of anchor state-action pairs $\{(s_i, a_i) : (s_i, a_i) \in \mathcal{K}\}$:*

$$\phi(s, a) = \sum_{i: (s_i, a_i) \in \mathcal{K}} \lambda_i(s, a) \phi(s_i, a_i) \quad \text{for} \quad \sum_{i=1}^K \lambda_i(s, a) = 1 \quad \text{and} \quad \lambda_i(s, a) \geq 0. \quad (2)$$

Further, we assume that the vectors in $\{\phi(s, a) : (s, a) \in \mathcal{K}\}$ are linearly independent.

We pause to develop some intuition of this assumption using Examples 1 and 2. In Example 1, it is straightforward to check that tabular MDPs satisfies Assumption 1 with $\mathcal{K} = \mathcal{S} \times \mathcal{A}$. In terms of Example 2, without loss of generality we can assume that the subspace spanned by the features has full rank, i.e. $\text{span}\{\phi(s, a) : (s, a) \in \mathcal{S} \times \mathcal{A}\} = \mathbb{R}^K$ (otherwise we can reduce the dimension of feature space). Then we can also check that Example 2 satisfies Assumption 1 with arbitrary $\mathcal{K} \subseteq \mathcal{S} \times \mathcal{A}$ such that the vectors in $\{\phi(s, a) : (s, a) \in \mathcal{K}\}$ are linearly independent. In fact, this sort of ‘‘anchor’’ notion appears widely in the literature: [3] considers ‘‘anchor word’’ in topic modeling; [19] defines ‘‘separability’’ in their study of non-negative matrix factorization; [48] introduces ‘‘aggregate’’ in reinforcement learning; [21] studies ‘‘anchor state’’ in soft state aggregation models. These concepts all bear some kind of resemblance to our definition of anchor state-action pairs here.

Throughout this paper, we assume that the feature mapping ϕ is known, which is a widely adopted assumption in previous literature [63, 30, 68, 26, 52, 56, 60]. In practice, large scale RL usually makes use of representation learning to obtain the feature mapping ϕ . Furthermore, the learned representations can be selected to satisfy the anchor state-action pairs assumption by design.

A useful implication of Assumption 1 is that we can represent the transition kernel as

$$P(\cdot|s, a) = \sum_{i: (s_i, a_i) \in \mathcal{K}} \lambda_i(s, a) P(\cdot|s_i, a_i), \quad (3)$$

This follows simply from substituting (2) into (1) (see (14) in Appendix A for a formal proof).

3 Model-based RL with a generative model

We start with studying model-based RL with a generative model in this section. We propose a model-based planning algorithm and show that it returns an ε -optimal policy with minimax optimal sample size.

3.1 Main results

A generative model and an empirical MDP. We assume access to a generative model that provides us with independent samples from M . For each anchor state-action pair $(s_i, a_i) \in \mathcal{K}$, we collect N

²Without loss of generality, one can always assume that the number of anchor state-action pairs equals to the feature dimension K . Interested readers are referred to Appendix D for detailed argument.

independent samples $s_i^{(j)} \sim P(\cdot | s_i, a_i)$, $j \in [N]$. This allows us to construct an empirical transition kernel \hat{P} where

$$\hat{P}(s' | s, a) = \sum_{i=1}^K \lambda_i(s, a) \cdot \left(\frac{1}{N} \sum_{j=1}^N \mathbb{1}\{s_i^{(j)} = s'\} \right), \quad (4)$$

for each $(s, a) \in \mathcal{S} \times \mathcal{A}$. Here, $\frac{1}{N} \sum_{j=1}^N \mathbb{1}\{s_i^{(j)} = s'\}$ is an empirical estimate of $P(s' | s_i, a_i)$ and then (3) is employed. With \hat{P} in hand, we can construct an empirical MDP $\hat{M} = (\mathcal{S}, \mathcal{A}, \hat{P}, r, \gamma)$. Our goal here is to derive the sample complexity which guarantees that the optimal policy of \hat{M} is an ε -optimal policy for the true MDP M . The algorithm is summarized below.

Algorithm 1 Model-based RL with a generative model

Inputs: a set of anchor state-action pairs $\mathcal{K} = \{(s_i, a_i) : i \in [K]\}$, feature mapping $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^K$, any planning algorithm \mathcal{P} , target algorithmic error level ε_{opt} .

For $i = 1, \dots, K$ **do**

Draw N independent samples $s_i^{(j)} \sim P(\cdot | s_i, a_i)$, $j = 1, \dots, N$.

End for

Construct an empirical MDP $\hat{M} = (\mathcal{S}, \mathcal{A}, \hat{P}, r, \gamma)$ where \hat{P} can be computed by (4).

Output $\hat{\pi}$ as an ε_{opt} -optimal policy of \hat{M} computed by the planning algorithm \mathcal{P} .

Careful readers may note that in Algorithm 1, $\{\lambda(s, a) : (s, a) \in \mathcal{S} \times \mathcal{A}\}$ is used in the construction of \hat{P} , while $\{\lambda(s, a) : (s, a) \in \mathcal{S} \times \mathcal{A}\}$ is not input into the algorithm. This is because given \mathcal{K} and ϕ are known, $\{\lambda(s, a) : (s, a) \in \mathcal{S} \times \mathcal{A}\}$ can be calculated explicitly. The following theorem provides theoretical guarantees for the output policy $\hat{\pi}$ of the chosen optimization algorithm on the empirical MDP \hat{M} .

Theorem 1. Suppose that $\delta > 0$ and $\varepsilon \in (0, (1-\gamma)^{-1/2}]$. Let $\hat{\pi}$ be the policy returned by Algorithm 1. Assume that

$$N \geq \frac{C \log(K / ((1-\gamma)\delta))}{(1-\gamma)^3 \varepsilon^2} \quad (5)$$

for some sufficiently large constant $C > 0$. Then with probability exceeding $1 - \delta$,

$$Q^*(s, a) - Q^{\hat{\pi}}(s, a) \leq \varepsilon + \frac{4\varepsilon_{\text{opt}}}{1-\gamma}, \quad (6)$$

for every $(s, a) \in \mathcal{S} \times \mathcal{A}$. Here ε_{opt} is the target algorithmic error level in Algorithm 1.

We first remark that the two terms on the right hand side of (6) can be viewed as statistical error and algorithmic error, respectively. The first term ε denotes the statistical error coming from the deviation of the empirical MDP \hat{M} from the true MDP M . As the sample size N grows, ε could decrease towards 0. The other term $4\varepsilon_{\text{opt}}/(1-\gamma)$ represents the algorithmic error where ε_{opt} is the target accuracy level of the planning algorithm applied to \hat{M} . Note that ε_{opt} can be arbitrarily small if we run the planning algorithm (e.g. value iteration) for enough iterations. A few implications of this theorem are in order.

- *Minimax-optimal sample complexity.* Assume that ε_{opt} is made negligibly small, e.g. $\varepsilon_{\text{opt}} = O((1-\gamma)\varepsilon)$ to be discussed in the next point. Note that we draw N independent samples for each state-action pair $(s, a) \in \mathcal{K}$, therefore the requirement (5) for finding an $O(\varepsilon)$ -optimal policy translates into the following sample complexity requirement

$$\tilde{O}\left(\frac{K}{(1-\gamma)^3 \varepsilon^2}\right).$$

This matches the minimax optimal lower bound (up to a logarithm factor) established in [63, Theorem 1] for feature-based MDP. In comparison, for tabular MDP the minimax optimal sample complexity is $\tilde{\Omega}((1-\gamma)^{-3} \varepsilon^{-2} |\mathcal{S}| |\mathcal{A}|)$ [5, 2]. Our sample complexity scales linearly with K instead of $|\mathcal{S}| |\mathcal{A}|$ for tabular MDP as desired.

- *Computational complexity.* An advantage of Theorem 1 is that it incorporates the use of any efficient planning algorithm applied to the empirical MDP \widehat{M} . Classical algorithms include Q-value iteration (QVI) or policy iteration (PI) [43]. For example, QVI achieves the target level ε_{opt} in $O((1-\gamma)^{-1} \log \varepsilon_{\text{opt}}^{-1})$ iterations, and each iteration takes time proportional to $O(NK + |\mathcal{S}||\mathcal{A}|K)$. To learn an $O(\varepsilon)$ -optimal policy, which requires sample complexity (5) and the target level $\varepsilon_{\text{opt}} = O((1-\gamma)\varepsilon)$, the overall running time is

$$\tilde{O} \left(\frac{|\mathcal{S}||\mathcal{A}|K}{1-\gamma} + \frac{K}{(1-\gamma)^4 \varepsilon^2} \right).$$

In comparison, for the tabular MDP the corresponding running time is $\tilde{O}((1-\gamma)^{-4} \varepsilon^{-2} |\mathcal{S}||\mathcal{A}|)$ [2]. This suggests that under the feature-based linear transition model, the computational complexity is $\min\{|\mathcal{S}||\mathcal{A}|/K, (1-\gamma)^{-3} \varepsilon^{-2}/K\}$ times lower than that for the tabular MDP (up to logarithm factors), which is significantly more efficient when K is not too large.

- *Stability vis-à-vis model misspecification.* A more general version of Theorem 1 (Theorem 3 in Appendix B) shows that when P approximately (instead of exactly) admits the linear transition model, we can still achieve some meaningful result. Specifically, if there exists a linear transition kernel \tilde{P} obeying $\max_{(s,a) \in \mathcal{S} \times \mathcal{A}} \|\tilde{P}(\cdot|s,a) - P(\cdot|s,a)\|_1 \leq \xi$ for some $\xi \geq 0$, we can show that $\hat{\pi}$ returned by Algorithm 1 (with slight modification) satisfies

$$Q^*(s,a) - Q^{\hat{\pi}}(s,a) \leq \varepsilon + \frac{4\varepsilon_{\text{opt}}}{1-\gamma} + \frac{22\xi}{(1-\gamma)^2},$$

for every $(s,a) \in \mathcal{S} \times \mathcal{A}$. This shows that the model-based method is stable vis-à-vis model misspecification. Interested readers are referred to Appendix B for more details.

In Algorithm 1, the reward function r is assumed to be known. If the information of r is unavailable, an alternative is to assume that r is linear with respect to the feature mapping ϕ , i.e. $r(s,a) = \theta^\top \phi(s,a)$ for every $(s,a) \in \mathcal{S} \times \mathcal{A}$, which is widely adopted in linear MDP literature [26, 30, 56, 60]. Under this linear assumption, one can obtain θ by solving the following linear system of equations

$$r(s,a) = \theta^\top \phi(s,a), \quad \forall (s,a) \in \mathcal{K}, \quad (7)$$

which can be constructed by the observed reward $r(s,a)$ for all anchor state-action pairs.

4 Model-free RL—vanilla Q Learning

In this section, we turn to study one of the most popular model-free RL algorithms—Q-learning. We provide tight sample complexity bound for vanilla Q-learning under the feature-based linear transition model, which shows its sample-efficiency (depends on $|K|$ instead of $|\mathcal{S}|$ or $|\mathcal{A}|$) and sub-optimality in the dependency on the effective horizon.

4.1 Q-learning algorithm

The vanilla Q-learning algorithm maintains a Q-function estimate $Q_t : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ for all $t \geq 0$, with initialization Q_0 obeying $0 \leq Q_0(s,a) \leq \frac{1}{1-\gamma}$ for every $(s,a) \in \mathcal{S} \times \mathcal{A}$. Assume we have access to a generative model. In each iteration $t \geq 1$, we collect an independent sample $s_t(s,a) \sim P(\cdot|s,a)$ for every anchor state-action pair $(s,a) \in \mathcal{K}$ and define $Q_{\mathcal{K}}^{(t)} : \mathcal{K} \rightarrow \mathbb{R}$ to be

$$Q_{\mathcal{K}}^{(t)}(s,a) := \max_{a' \in \mathcal{A}} Q_t(s_t, a'), \quad s_t \equiv s_t(s,a) \sim P(\cdot|s,a).$$

Then given the learning rate $\eta_t \in (0, 1]$, the algorithm adopts the following update rule to update all entries of the Q-function estimate

$$Q_t = (1 - \eta_t) Q_{t-1} + \eta_t \mathcal{T}_{\mathcal{K}}^{(t)}(Q_{t-1}).$$

Here, $\mathcal{T}_{\mathcal{K}}^{(t)}$ is an empirical Bellman operator associated with the linear transition model M and the set \mathcal{K} and is given by

$$\mathcal{T}_{\mathcal{K}}^{(t)}(Q)(s,a) := r(s,a) + \gamma \lambda(s,a) Q_{\mathcal{K}}^{(t)},$$

where (3) is used in the construction. Clearly, this newly defined operator $\mathcal{T}_{\mathcal{K}}^{(t)}$ is an unbiased estimate of the famous Bellman operator \mathcal{T} [8] defined as

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \mathcal{T}(Q)(s, a) := r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\max_{a' \in \mathcal{A}} Q(s', a') \right].$$

A critical property is that the Bellman operator \mathcal{T} is contractive with a unique fixed point which is the optimal Q-function Q^* [8]. To solve the fixed-point equation $\mathcal{T}(Q^*) = Q^*$, Q-learning was then introduced by [58] based on the idea of stochastic approximation [44]. This procedure is precisely described in Algorithm 2.

Algorithm 2 Vanilla Q-learning for infinite-horizon discounted MDPs

inputs: learning rates $\{\eta_t\}$, number of iterations T , discount factor γ , initial estimate Q_0 .

for $t = 1, \dots, T$ **do**

 Draw $s_t(s, a) \sim P(\cdot|s, a)$ for each $(s, a) \in \mathcal{K}$.

 Compute Q_t according to the update rule

$$Q_t = (1 - \eta_t) Q_{t-1} + \eta_t \mathcal{T}_{\mathcal{K}}^{(t)}(Q_{t-1}).$$

end for

4.2 Main results

We are now ready to provide our main result for vanilla Q-learning, assuming sampling access to a generative model.

Theorem 2. *Consider any $\delta \in (0, 1)$ and $\varepsilon \in (0, 1]$. Assume that for any $0 \leq t \leq T$, the learning rates satisfy*

$$\frac{1}{1 + \frac{c_1(1-\gamma)T}{\log^2 T}} \leq \eta_t \leq \frac{1}{1 + \frac{c_2(1-\gamma)t}{\log^2 T}} \quad (8)$$

for some sufficiently small universal constants $c_1 \geq c_2 > 0$. Suppose that the total number of iterations T exceeds

$$T \geq \frac{C_3 \log(KT/\delta) \log^4 T}{(1-\gamma)^4 \varepsilon^2} \quad (9)$$

for some sufficiently large universal constant $C_3 > 0$. If the initialization obeys $0 \leq Q_0(s, a) \leq \frac{1}{1-\gamma}$ for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, then with probability exceeding $1 - \delta$, the output Q_T of Algorithm 2 satisfies

$$\max_{(s, a) \in \mathcal{S} \times \mathcal{A}} |Q_T(s, a) - Q^*(s, a)| \leq \varepsilon. \quad (10)$$

In addition, let π_T (resp. V_T) to be the policy (resp. value function) induced by Q_T , then one has

$$\max_{s \in \mathcal{S}} |V^{\pi_T}(s) - V^*(s)| \leq \frac{2\gamma\varepsilon}{1-\gamma}. \quad (11)$$

This theorem provides theoretical guarantees on the performance of Algorithm 2. A few implications of this theorem are in order.

- *Learning rate.* The condition (8) accommodates two commonly adopted choice of learning rates: (i) linearly rescaled learning rates $\eta_t = [1 + c_2(1-\gamma)t/\log^2 T]^{-1}$, and (ii) iteration-invariant learning rates $\eta_t \equiv [1 + c_1(1-\gamma)T/\log^2 T]$. Interested readers are referred to the discussions in [34, Section 3.1] for more details on these two learning rate schemes.
- *Tight sample complexity bound.* Note that we draw K independent samples in each iteration, therefore the iteration complexity (9) can be translated into the sample complexity bound TK in order for Q-learning to achieve ε -accuracy:

$$\tilde{O} \left(\frac{K}{(1-\gamma)^4 \varepsilon^2} \right). \quad (12)$$

As we will see shortly, this result improves the state-of-the-art sample complexity bound presented in [63, Theorem 2]. In addition, the dependency on the effective horizon $(1-\gamma)^{-4}$ matches the lower bound established in [34, Theorem 2] for vanilla Q-learning using either learning rate scheme covered in the previous remark, suggesting that our sample complexity bound (12) is sharp.

- *Stability vis-à-vis model misspecification.* Just like the model-based approach, we can also show that Q-learning is also stable vis-à-vis model misspecification when P approximately admits the linear transition model. We refer interested readers to Theorem 4 in Appendix B for more details.

Comparison with [63]. We compare our result with the sample complexity bounds for Q-learning under the feature-based linear transition model in [63].

- We first compare our result with [63, Theorem 2], which is, to the best of our knowledge, the state-of-the-art theory for this problem. When there is no model misspecification, [63, Theorem 2] showed that in order for their Phased Parametric Q-learning³ (Algorithm 1 therein) to learn an ε -optimal policy, the sample size needs to be

$$\tilde{O}\left(\frac{K}{(1-\gamma)^7 \varepsilon^2}\right).$$

Note that (12) is the sample complexity required for entrywise ε -accurate estimate of the optimal Q-function, thus a fair comparison requires to use the sample complexity for learning an ε -optimal policy deduced from (11), which is

$$\tilde{O}\left(\frac{K}{(1-\gamma)^6 \varepsilon^2}\right).$$

Hence, our sample complexity improves upon previous work by a factor at least on the order of $(1-\gamma)^{-1}$. However it is worth mentioning that [63, Theorem 2] is built upon weaker conditions $\sum_{i=1}^K \lambda_i(s, a) = 1$ and $\sum_{i=1}^K |\lambda_i(s, a)| \leq L$ for some $L \geq 1$, which does not require $\lambda_i(s, a) \geq 0$. Our result holds under Assumption 1, which requires $\sum_{i=1}^K \lambda_i(s, a) = 1$ and $\lambda_i(s, a) \geq 0$. Under the current analysis framework, it is difficult to obtain tight sample complexity bounds without assuming $\lambda_i(s, a) \geq 0$.

- Besides vanilla Q-learning, [63] also proposed a new variant of Q-learning called Optimal Phased Parametric Q-Learning (Algorithm 2 therein), which is essentially Q-learning with variance reduction. [63, Theorem 3] showed that the sample complexity for this algorithm is

$$\tilde{O}\left(\frac{K}{(1-\gamma)^3 \varepsilon^2}\right),$$

which matches minimax optimal lower bound (up to a logarithm factor) established in [63, Theorem 1]. Careful reader might notice that this sample complexity bound is better than ours for vanilla Q-learning. We emphasize that as elucidated in the second implication under Theorem 2, our result is already tight for vanilla Q-learning. This observation reveals that while the sample complexity for vanilla Q-learning is provably sub-optimal, the variants of Q-learning can have better performance and achieve minimax optimal sample complexity.

We conclude this section by comparing model-based and model-free approaches. Theorem 1 shows that the sample complexity of the model-based approach is minimax optimal, whilst vanilla Q-learning, perhaps the most commonly adopted model-free method, is sub-optimal according to Theorem 2. However this does not mean that model-based method is better than model-free ones since (i) some variants of Q-learning (see [63, Algorithm 2] for example) also has minimax optimal sample complexity; and (ii) in many applications it might be unrealistic to estimate the model in advance.

5 A glimpse of our technical approaches

The establishment of Theorems 1 and 2 calls for a series of technical novelties in the proof. In what follows, we briefly highlight our key technical ideas and novelties.

³The difference between Algorithm 2 and Phased Parametric Q-Learning in [63] is that Algorithm 2 maintains and updates a Q-function estimate Q_t , while Phased Parametric Q-Learning parameterized Q-function by

$$Q_w(s, a) := r(s, a) + \gamma \phi(s, a)^\top w,$$

and then updates the parameters w .

- For the model-based approach, we employ “leave-one-out” analysis to decouple the complicated statistical dependency between the empirical probability transition model \hat{P} and the corresponding optimal policy. Specifically, [2] proposed to construct a collection of auxiliary MDPs where each one of them leaves out a single state s by setting s to be an absorbing state and keeping everything else untouched. We tailor this high level idea to the needs of linear transition model, then the independence between the newly constructed MDP with absorbing state s and data samples collected at state s will facilitate our analysis, as detailed in Lemma 1. Compared with [2], Theorem 1 extends the tabular setting studied in [2] to the linear transition model and accommodates model misspecification, which actually needs significant efforts as detailed in the supplementary materials. This “leave-one-out” type of analysis has been utilized in studying numerous problems by a long line of work, such as [22, 38, 53, 13, 12, 14, 15], just to name a few.
- To obtain tighter sample complexity bound than the previous one $\tilde{O}(\frac{K}{(1-\gamma)^7 \varepsilon^2})$ in [63] for vanilla Q-learning, we invoke Freedman’s inequality [24] for the concentration of an error term with martingale structure as illustrated in Appendix C, while the classical ones used in analyzing Q-learning are Hoeffding’s inequality and Bernstein’s inequality [63]. The use of Freedman’s inequality helps us establish a recursive relation on $\{\|Q_t - Q^*\|_\infty\}_{t=0}^T$, which consequently leads to the performance guarantee (10). It is worth mentioning that [34] also studied vanilla Q-learning in the tabular MDP setting and adopted Freedman’s inequality, while we emphasize that it requires a lot of efforts and more delicate analyses in order to study linear transition model and also allow for model misspecification in the current paper, as detailed in the supplementary material.

6 Additional related literature

To remedy the issue of prohibitively high sample complexity, there exists a substantial body of literature proposing and studying many structural assumptions and complexity notions under different settings. This current paper focuses on linear transition model which is studied in MDP by numerous previous works [63, 30, 64, 68, 40, 25, 56, 52, 26, 60]. Among them, [63] studied linear transition model and provided tight sample complexity bounds for a new variant of Q-learning with the help of variance reduction. [30] focused on linear MDP and designed an algorithm called “Least-Squares Value Iteration with UCB” with both polynomial runtime and polynomial sample complexity without accessing generative model. [56] extended the study of linear MDP to the framework of reward-free reinforcement learning. [68] considered a different feature mapping called linear kernel MDP and devised an algorithm with polynomial regret bound without generative model. Other popular structure assumptions include: [61] studied fully deterministic transition dynamics; [28] introduced Bellman rank and proposed an algorithm which needs sample size polynomially dependent on Bellman rank to obtain a near-optimal policy in contextual decision processes; [20] assumed that the value function has low variance compared to the mean for all deterministic policy; [39, 42, 7, 66] used linear model to approximate the value function; [35] assumed that the optimal Q-function can be linearly-parameterized by the features.

Apart from the linear transition model, another notion adopted in this work is the generative model, whose role in discounted MDP has been studied by extensive literature. The concept of generative model was originally introduced by [32], and then widely adopted in numerous works, including [31, 5, 63, 54, 2, 41], to name a few. Specifically, it is assumed that a generative model of the studied MDP is available and can be queried for every state-action pair and output the next state. Among previous works, [5] proved that the minimax lower bound on the sample complexity to obtain an ε -optimal policy was $\tilde{\Omega}(\frac{|S||A|}{(1-\gamma)^3 \varepsilon^2})$. [5] also showed that model-based approach can output an ε -optimal value function with near-optimal sample complexity for $\varepsilon \in (0, 1)$. Then [2] made significant progress on the challenging problem of establishing minimax optimal sample complexity in estimating an ε -optimal policy with the help of “leave-one-out” analysis.

In addition, after being proposed in [59], Q-learning has become the focus of a rich line of research [58, 11, 32, 23, 4, 29, 53, 16, 37, 62]. Among them, [16, 37, 62] studied Q-learning in the presence of Markovian data, i.e. a single sample trajectory. In contrast, under the generative setting of Q-learning where a fresh sample can be drawn from the simulator at each iteration, [54] analyzed a variant of Q-learning with the help of variance reduction, which was proved to enjoy minimax optimal sample complexity $\tilde{O}(\frac{|S||A|}{(1-\gamma)^3 \varepsilon^2})$. Then more recently, [34] improved the lower bound of the

vanilla Q-learning algorithm in terms of its scaling with $\frac{1}{1-\gamma}$ and proved a matching upper bound $\tilde{O}(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^4 \varepsilon^2})$.

7 Discussion

This paper studies sample complexity of both model-based and model-free RL under a discounted infinite-horizon MDP with feature-based linear transition model. We establish tight sample complexity bounds for both model-based approaches and Q-learning, which scale linearly with the feature dimension K instead of $|\mathcal{S}| \times |\mathcal{A}|$, thus considerably reduce the required sample size for large-scale MDPs when K is relatively small. Our results are sharp, and the sample complexity bound for the model-based approach matches the minimax lower bound. The current work suggests a couple of directions for future investigation, as discussed in detail below.

- *Extension to episodic MDPs.* An interesting direction for future research is to study linear transition model in episodic MDP. This focus of this work is infinite-horizon discounted MDPs, and hopefully the analysis here can be extended to study the episodic MDP as well ([17, 18, 6, 27, 55, 26]).
- *Continuous state and action space.* The state and action spaces in this current paper are still assumed to be finite, since the proof relies heavily on the matrix operations. However, we expect that the results can be extended to accommodate continuous state and action space by employing more complicated analysis.
- *Accommodating entire range of ε .* Since both value functions and Q-functions can take value in $[0, (1-\gamma)^{-1}]$, ideally our theory should cover all choices of $\varepsilon \in (0, (1-\gamma)^{-1}]$. However we require that $\varepsilon \in (0, (1-\gamma)^{-1/2}]$ in Theorem 1 and $\varepsilon \in (0, 1]$ in Theorem 2. While most of the prior works like [2, 63] also impose these restrictions, a recent work [36] proposed a perturbed model-based planning algorithm and proved minimax optimal guarantees for any $\varepsilon \in (0, (1-\gamma)^{-1}]$. While their work only focused on model-based RL under tabular MDP, an interesting future direction is to improve our theory to accommodate any $\varepsilon \in (0, (1-\gamma)^{-1}]$.
- *General function approximation.* Another future direction is to extend the study to more general function approximation starting from linear structure covered in this paper. There exists a rich body of work proposing and studying different structures, such as linear value function approximation [39, 42, 7, 66], linear MDPs with infinite dimensional features [1], Eluder dimension [57], Bellman rank [28] and Witness rank [50], etc. Therefore, it is hopeful to investigate these settings and improve the sample efficiency.

Acknowledgments and Disclosure of Funding

B. Wang is supported in part by Gordon Y. S. Wu Fellowships in Engineering. Y. Yan is supported in part by ARO grant W911NF-20-1-0097 and NSF grant CCF-1907661. Part of this work was done while Y. Yan was visiting the Simons Institute for the Theory of Computing. J. Fan is supported in part by the ONR grant N00014-19-1-2120 and the NSF grants DMS-1662139, DMS-1712591, DMS-2052926, DMS-2053832, and the NIH grant 2R01-GM072611-15.

References

- [1] Alekh Agarwal, Mikael Henaff, Sham Kakade, and Wen Sun. Pc-pg: Policy cover directed exploration for provable policy gradient learning. *arXiv preprint arXiv:2007.08459*, 2020.
- [2] Alekh Agarwal, Sham Kakade, and Lin F Yang. Model-based reinforcement learning with a generative model is minimax optimal. In *Conference on Learning Theory*, pages 67–83. PMLR, 2020.
- [3] Sanjeev Arora, Rong Ge, and Ankur Moitra. Learning topic models—going beyond svd. In *2012 IEEE 53rd annual symposium on foundations of computer science*, pages 1–10. IEEE, 2012.
- [4] Mohammad Gheshlaghi Azar, Remi Munos, M Ghavamzadeh, and Hilbert J Kappen. Speedy q-learning. 2011.

- [5] Mohammad Gheshlaghi Azar, Rémi Munos, and Hilbert J Kappen. Minimax pac bounds on the sample complexity of reinforcement learning with a generative model. *Machine learning*, 91(3):325–349, 2013.
- [6] Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In *International Conference on Machine Learning*, pages 263–272. PMLR, 2017.
- [7] Kamyar Azizzadenesheli, Emma Brunskill, and Animashree Anandkumar. Efficient exploration through bayesian deep q-networks. In *2018 Information Theory and Applications Workshop (ITA)*, pages 1–9. IEEE, 2018.
- [8] Richard Bellman. On the theory of dynamic programming. *Proceedings of the National Academy of Sciences of the United States of America*, 38(8):716, 1952.
- [9] Richard Bellman and Stuart Dreyfus. Functional approximations and dynamic programming. *Mathematical Tables and Other Aids to Computation*, pages 247–251, 1959.
- [10] Dimitri P Bertsekas et al. *Dynamic programming and optimal control: Vol. 1*. Athena scientific Belmont, 2000.
- [11] Dimitri P Bertsekas and John N Tsitsiklis. *Neuro-dynamic programming*. Athena Scientific, 1996.
- [12] Yuxin Chen, Yuejie Chi, Jianqing Fan, Cong Ma, and Yuling Yan. Noisy matrix completion: Understanding statistical guarantees for convex relaxation via nonconvex optimization. *SIAM Journal on Optimization*, 30(4):3098–3121, 2020.
- [13] Yuxin Chen, Jianqing Fan, Cong Ma, and Kaizheng Wang. Spectral method and regularized mle are both optimal for top-k ranking. *Annals of statistics*, 47(4):2204, 2019.
- [14] Yuxin Chen, Jianqing Fan, Cong Ma, and Yuling Yan. Bridging convex and nonconvex optimization in robust pca: Noise, outliers, and missing data. *arXiv preprint arXiv:2001.05484*, accepted to *Annals of Statistics*, 2020.
- [15] Yuxin Chen, Jianqing Fan, Bingyan Wang, and Yuling Yan. Convex and nonconvex optimization are both minimax-optimal for noisy blind deconvolution under random designs. *Journal of the American Statistical Association*, (just-accepted):1–27, 2021.
- [16] Zaiwei Chen, Sheng Zhang, Thinh T Doan, Siva Theja Maguluri, and John-Paul Clarke. Performance of q-learning with linear function approximation: Stability and finite-time analysis. *arXiv preprint arXiv:1905.11425*, 2019.
- [17] Christoph Dann and Emma Brunskill. Sample complexity of episodic fixed-horizon reinforcement learning. *Advances in Neural Information Processing Systems*, 28:2818–2826, 2015.
- [18] Christoph Dann, Tor Lattimore, and Emma Brunskill. Unifying pac and regret: uniform pac bounds for episodic reinforcement learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 5717–5727, 2017.
- [19] David Donoho and Victoria Stodden. When does non-negative matrix factorization give a correct decomposition into parts? In *17th Annual Conference on Neural Information Processing Systems, NIPS 2003*. Neural information processing systems foundation, 2004.
- [20] Simon S Du, Yuping Luo, Ruosong Wang, and Hanrui Zhang. Provably efficient q-learning with function approximation via distribution shift error checking oracle. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pages 8060–8070, 2019.
- [21] Yaqi Duan, Zheng Tracy Ke, and Mengdi Wang. State aggregation learning from markov transition data. *Advances in Neural Information Processing Systems*, 32, 2019.
- [22] Noureddine El Karoui. On the impact of predictor geometry on the performance on high-dimensional ridge-regularized generalized robust regression estimators. *Probability Theory and Related Fields*, 170(1):95–175, 2018.

- [23] Eyal Even-Dar, Yishay Mansour, and Peter Bartlett. Learning rates for q-learning. *Journal of machine learning Research*, 5(1), 2003.
- [24] David A Freedman. On tail probabilities for martingales. *the Annals of Probability*, pages 100–118, 1975.
- [25] Botao Hao, Yaqi Duan, Tor Lattimore, Csaba Szepesvári, and Mengdi Wang. Sparse feature selection makes batch reinforcement learning more sample efficient. In *International Conference on Machine Learning*, pages 4063–4073. PMLR, 2021.
- [26] Jiafan He, Dongruo Zhou, and Quanquan Gu. Logarithmic regret for reinforcement learning with linear function approximation. In *International Conference on Machine Learning*, pages 4171–4180. PMLR, 2021.
- [27] Nan Jiang and Alekh Agarwal. Open problem: The dependence of sample complexity lower bounds on planning horizon. In *Conference On Learning Theory*, pages 3395–3398. PMLR, 2018.
- [28] Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Contextual decision processes with low bellman rank are pac-learnable. In *International Conference on Machine Learning*, pages 1704–1713. PMLR, 2017.
- [29] Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is q-learning provably efficient? In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 4868–4878, 2018.
- [30] Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143. PMLR, 2020.
- [31] Sham Machandranath Kakade. *On the sample complexity of reinforcement learning*. PhD thesis, UCL (University College London), 2003.
- [32] Michael Kearns and Satinder Singh. Finite-sample convergence rates for q-learning and indirect algorithms. *Advances in neural information processing systems*, pages 996–1002, 1999.
- [33] B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 2021.
- [34] Gen Li, Changxiao Cai, Yuxin Chen, Yuantao Gu, Yuting Wei, and Yuejie Chi. Is q-learning minimax optimal? a tight sample complexity analysis. *arXiv preprint arXiv:2102.06548*, 2021.
- [35] Gen Li, Yuxin Chen, Yuejie Chi, Yuantao Gu, and Yuting Wei. Sample-efficient reinforcement learning is feasible for linearly realizable mdps with limited revisiting. *arXiv preprint arXiv:2105.08024*, 2021.
- [36] Gen Li, Yuting Wei, Yuejie Chi, Yuantao Gu, and Yuxin Chen. Breaking the sample size barrier in model-based reinforcement learning with a generative model. *Advances in Neural Information Processing Systems*, 33, 2020.
- [37] Gen Li, Yuting Wei, Yuejie Chi, Yuantao Gu, and Yuxin Chen. Sample complexity of asynchronous q-learning: Sharper analysis and variance reduction. *Advances in neural information processing systems*, 2020.
- [38] Cong Ma, Kaizheng Wang, Yuejie Chi, and Yuxin Chen. Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval and matrix completion. In *International Conference on Machine Learning*, pages 3345–3354. PMLR, 2018.
- [39] Francisco S Melo and M Isabel Ribeiro. Q-learning with linear function approximation. In *International Conference on Computational Learning Theory*, pages 308–322. Springer, 2007.

- [40] Aditya Modi, Nan Jiang, Ambuj Tewari, and Satinder Singh. Sample complexity of reinforcement learning using linearly combined model ensembles. In *International Conference on Artificial Intelligence and Statistics*, pages 2010–2020. PMLR, 2020.
- [41] Ashwin Pananjady and Martin J Wainwright. Instance-dependent ℓ_∞ -bounds for policy evaluation in tabular reinforcement learning. *IEEE Transactions on Information Theory*, 67(1):566–585, 2020.
- [42] Ronald Parr, Lihong Li, Gavin Taylor, Christopher Painter-Wakefield, and Michael L Littman. An analysis of linear models, linear value-function approximation, and feature selection for reinforcement learning. In *Proceedings of the 25th international conference on Machine learning*, pages 752–759, 2008.
- [43] Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- [44] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- [45] Roshan Shariff and Csaba Szepesvári. Efficient planning in large mdps with weak linear function approximation. *arXiv preprint arXiv:2007.06184*, 2020.
- [46] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- [47] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- [48] Satinder P Singh, Tommi Jaakkola, and Michael I Jordan. Reinforcement learning with soft state aggregation. *Advances in neural information processing systems* 7, 7:361, 1995.
- [49] Satinder P Singh and Richard C Yee. An upper bound on the loss from approximate optimal-value functions. *Machine Learning*, 16(3):227–233, 1994.
- [50] Wen Sun, Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Model-based rl in contextual decision processes: Pac bounds and exponential improvements over model-free approaches. In *Conference on Learning Theory*, pages 2898–2933. PMLR, 2019.
- [51] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [52] Ahmed Touati and Pascal Vincent. Efficient learning in non-stationary linear markov decision processes. *arXiv preprint arXiv:2010.12870*, 2020.
- [53] Martin J Wainwright. Stochastic approximation with cone-contractive operators: Sharp ℓ_∞ -bounds for q -learning. *arXiv preprint arXiv:1905.06265*, 2019.
- [54] Martin J Wainwright. Variance-reduced q -learning is minimax optimal. *arXiv preprint arXiv:1906.04697*, 2019.
- [55] Ruosong Wang, Simon S Du, Lin Yang, and Sham Kakade. Is long horizon rl more difficult than short horizon rl? *Advances in Neural Information Processing Systems*, 33, 2020.
- [56] Ruosong Wang, Simon S Du, Lin F Yang, and Ruslan Salakhutdinov. On reward-free reinforcement learning with linear function approximation. *arXiv preprint arXiv:2006.11274*, 2020.
- [57] Ruosong Wang, Russ R Salakhutdinov, and Lin Yang. Reinforcement learning with general value function approximation: Provably efficient approach via bounded eluder dimension. *Advances in Neural Information Processing Systems*, 33, 2020.

- [58] Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3-4):279–292, 1992.
- [59] Christopher John Cornish Hellaby Watkins. Learning from delayed rewards. 1989.
- [60] Chen-Yu Wei, Mehdi Jafarnia Jahromi, Haipeng Luo, and Rahul Jain. Learning infinite-horizon average-reward mdps with linear function approximation. In *International Conference on Artificial Intelligence and Statistics*, pages 3007–3015. PMLR, 2021.
- [61] Zheng Wen and Benjamin Van Roy. Efficient reinforcement learning in deterministic systems with value function generalization. *Mathematics of Operations Research*, 42(3):762–782, 2017.
- [62] Pan Xu and Quanquan Gu. A finite-time analysis of q-learning with neural network function approximation. In *International Conference on Machine Learning*, pages 10555–10565. PMLR, 2020.
- [63] Lin Yang and Mengdi Wang. Sample-optimal parametric q-learning using linearly additive features. In *International Conference on Machine Learning*, pages 6995–7004. PMLR, 2019.
- [64] Lin Yang and Mengdi Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. In *International Conference on Machine Learning*, pages 10746–10756. PMLR, 2020.
- [65] Andrea Zanette, David Brandfonbrener, Emma Brunskill, Matteo Pirodda, and Alessandro Lazaric. Frequentist regret bounds for randomized least-squares value iteration. In *International Conference on Artificial Intelligence and Statistics*, pages 1954–1964. PMLR, 2020.
- [66] Andrea Zanette, Alessandro Lazaric, Mykel Kochenderfer, and Emma Brunskill. Learning near optimal policies with low inherent bellman error. In *International Conference on Machine Learning*, pages 10978–10989. PMLR, 2020.
- [67] Andrea Zanette, Alessandro Lazaric, Mykel J Kochenderfer, and Emma Brunskill. Limiting extrapolation in linear approximate value iteration. *Advances in Neural Information Processing Systems*, 32:5615–5624, 2019.
- [68] Dongruo Zhou, Jiafan He, and Quanquan Gu. Provably efficient reinforcement learning for discounted mdps with feature mapping. In *International Conference on Machine Learning*, pages 12793–12802. PMLR, 2021.

Supplement to “Sample-Efficient Reinforcement Learning for Linearly-Parameterized MDPs with a Generative Model”

Bingyan Wang*
Princeton University
bingyanw@princeton.edu

Yuling Yan*
Princeton University
yulingy@princeton.edu

Jianqing Fan
Princeton University
jqfan@princeton.edu

A Notations

In this section we gather the notations that will be used throughout the appendix.

For any vectors $\mathbf{u} = [u_i]_{i=1}^n \in \mathbb{R}^n$ and $\mathbf{v} = [v_i]_{i=1}^n \in \mathbb{R}^n$, let $\mathbf{u} \circ \mathbf{v} = [u_i v_i]_{i=1}^n$ denote the Hadamard product of \mathbf{u} and \mathbf{v} . We slightly abuse notations to use $\sqrt{\cdot}$ and $|\cdot|$ to define entry-wise operation, i.e. for any vector $\mathbf{v} = [v_i]_{i=1}^n$ denote $\sqrt{\mathbf{v}} := [\sqrt{v_i}]_{i=1}^n$ and $|\mathbf{v}| := [|v_i|]_{i=1}^n$. Furthermore, the binary notations \leq and \geq are both defined in entry-wise manner, i.e. $\mathbf{u} \leq \mathbf{v}$ (resp. $\mathbf{u} \geq \mathbf{v}$) means $u_i \leq v_i$ (resp. $u_i \geq v_i$) for all $1 \leq i \leq n$. For a collection of vectors $\mathbf{v}_1, \dots, \mathbf{v}_m \in \mathbb{R}^n$ with $\mathbf{v}_i = [v_{i,j}]_{j=1}^n \in \mathbb{R}^n$, we define the max operator to be $\max_{1 \leq i \leq m} \mathbf{v}_i := [\max_{1 \leq i \leq m} v_{i,j}]_{j=1}^n$.

For any matrix $\mathbf{M} \in \mathbb{R}^{m \times n}$, $\|\mathbf{M}\|_1$ is defined as the largest row-wise ℓ_1 norm of \mathbf{M} , i.e. $\|\mathbf{M}\|_1 := \max_i \sum_j |M_{i,j}|$. In addition, we define $\mathbf{1}$ to be a vector with all the entries being 1, and \mathbf{I} be the identity matrix. To express the probability transition function P in matrix form, we define the matrix $\mathbf{P} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$ to be a matrix whose (s, a) -th row $\mathbf{P}_{s,a}$ corresponds to $P(\cdot|s, a)$. In addition, we define \mathbf{P}^π to be the probability transition matrix induced by policy π , i.e. $\mathbf{P}_{(s,a),(s',a')}^\pi = \mathbf{P}_{s,a}(s') \mathbb{1}_{\pi(s')=a'}$ for all state-action pairs (s, a) and (s', a') . We define π_t to be the policy induced by Q_t , i.e. $Q_t(s, \pi_t(s)) = \max_a Q_t(s, a)$ for all $s \in \mathcal{S}$. Furthermore, we denote the reward function r by vector $\mathbf{r} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$, i.e. the (s, a) -th element of \mathbf{r} equals $r(s, a)$. In the same manner, we define $\mathbf{V}^\pi \in \mathbb{R}^{|\mathcal{S}|}$, $\mathbf{V}^* \in \mathbb{R}^{|\mathcal{S}|}$, $\mathbf{V}_t \in \mathbb{R}^{|\mathcal{S}|}$, $\mathbf{Q}^\pi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$, $\mathbf{Q}^* \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ and $\mathbf{Q}_t \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ to represent V^π , V^* , V_t , Q^π , Q^* and Q_t respectively. By using these notations, we can rewrite the Bellman equation as

$$\mathbf{Q}^\pi = \mathbf{r} + \gamma \mathbf{P} \mathbf{V}^\pi = \mathbf{r} + \gamma \mathbf{P}^\pi \mathbf{Q}^\pi. \quad (11)$$

Further, for any vector $\mathbf{V} \in \mathbb{R}^{|\mathcal{S}|}$, let $\text{Var}_{\mathbf{P}}(\mathbf{V}) \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ be

$$\text{Var}_{\mathbf{P}}(\mathbf{V}) := \mathbf{P}(\mathbf{V} \circ \mathbf{V}) - (\mathbf{P}\mathbf{V}) \circ (\mathbf{P}\mathbf{V}), \quad (12)$$

and define $\text{Var}_{\mathbf{P}_{s,a}}(\mathbf{V}) \in \mathbb{R}$ to be

$$\text{Var}_{\mathbf{P}_{s,a}}(\mathbf{V}) := \mathbf{P}_{s,a}(\mathbf{V} \circ \mathbf{V}) - (\mathbf{P}_{s,a}\mathbf{V})^2, \quad (13)$$

where $\mathbf{P}_{s,a}$ is the (s, a) -th row of \mathbf{P} .

Next, we reconsider Assumption 1. For any state-action pair (s, a) , we define vector $\boldsymbol{\lambda}(s, a) \in \mathbb{R}^K$ (resp. $\boldsymbol{\phi}(s, a) \in \mathbb{R}^K$) with $\boldsymbol{\lambda}(s, a) = [\lambda_i(s, a)]_{i=1}^K$ (resp. $\boldsymbol{\phi}(s, a) = [\phi_i(s, a)]_{i=1}^K$) and matrix

*Equal contribution.

$\Lambda \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}| \times K}$ (resp. $\Phi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}| \times K}$) whose (s, a) -th row corresponds to $\lambda(s, a)^\top$ (resp. $\phi(s, a)^\top$). Define vector $\psi(s, a) \in \mathbb{R}^K$ with $\psi(s, a) = [\psi_i(s, a)]_{i=1}^K$ and matrix $\Psi \in \mathbb{R}^{K \times |\mathcal{S}|}$ whose (s, a) -th column corresponds to $\psi(s, a)^\top$. Further, let $P_K \in \mathbb{R}^{K \times |\mathcal{S}|}$ (resp. $\Phi_K \in \mathbb{R}^{K \times K}$) to be a submatrix of P (resp. Φ) formed by concatenating the rows $\{P_{s,a}, (s, a) \in \mathcal{K}\}$ (resp. $\{\Phi_{s,a}, (s, a) \in \mathcal{K}\}$). By using the previous notations, we can express the relations in Definition 1 and Assumption 1 as $P_K = \Phi_K \Psi$, $P = \Phi \Psi$ and $\Phi = \Lambda \Phi_K$. Note that Assumption 1 suggests Φ_K is invertible. Taking these equations collectively yields

$$P = \Phi \Psi = \Phi \Phi_K^{-1} P_K = \Lambda \Phi_K \Phi_K^{-1} P_K = \Lambda P_K, \quad (14)$$

which is reminiscent of the anchor word condition in topic modelling [2]. In addition, for each iteration t , we denote the collected samples as $\{s_t(s, a)\}_{(s,a) \in \mathcal{K}}$ and define a matrix $\hat{P}_K^{(t)} \in \{0, 1\}^{K \times |\mathcal{S}|}$ to be

$$\hat{P}_K^{(t)}((s, a), s') := \begin{cases} 1, & \text{if } s' = s_t(s, a) \\ 0, & \text{otherwise} \end{cases} \quad (15)$$

for any $(s, a) \in \mathcal{K}$ and $s' \in \mathcal{S}$. Further, we define $\hat{P}_t = \Lambda \hat{P}_K^{(t)}$. Then it is obvious to see that \hat{P}_t has nonnegative entries and unit ℓ_1 norm for each row due to Assumption 1, i.e. $\|\hat{P}_t\|_1 = 1$.

B Analysis of model-based RL (Proof of Theorem 1)

In this section, we will provide complete proof for Theorem 1. As a matter of fact, our proof strategy here justifies a more general version of Theorem 1 that accounts for model misspecification, as stated below.

Theorem 3. Suppose that $\delta > 0$ and $\varepsilon \in (0, (1 - \gamma)^{-1/2}]$. Assume that there exists a probability transition model \tilde{P} obeying Definition 1 and Assumption 1 with feature vectors $\{\phi(s, a)\}_{(s,a) \in \mathcal{S} \times \mathcal{A}} \subset \mathbb{R}^K$ and anchor state-action pairs \mathcal{K} such that

$$\|\tilde{P} - P\|_1 \leq \xi$$

for some $\xi \geq 0$. Let $\hat{\pi}$ be the policy returned by Algorithm 1. Assume that

$$N \geq \frac{C \log(K / ((1 - \gamma)\delta))}{(1 - \gamma)^3 \varepsilon^2} \quad (16)$$

for some sufficiently large constant $C > 0$. Then with probability exceeding $1 - \delta$,

$$Q^*(s, a) - Q^{\hat{\pi}}(s, a) \leq \varepsilon + \frac{4\varepsilon_{\text{opt}}}{1 - \gamma} + \frac{22\xi}{(1 - \gamma)^2}, \quad (17)$$

for every state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$.

Theorem 3 subsumes Theorem 1 as a special case with $\xi = 0$. The remainder of this section is devoted to proving Theorem 3.

B.1 Proof of Theorem 3

The error $Q^{\hat{\pi}} - Q^*$ can be decomposed as

$$\begin{aligned} Q^{\hat{\pi}} - Q^* &= Q^{\hat{\pi}} - \hat{Q}^{\hat{\pi}} + \hat{Q}^{\hat{\pi}} - \hat{Q}^* + \hat{Q}^* - Q^* \\ &\geq Q^{\hat{\pi}} - \hat{Q}^{\hat{\pi}} + \hat{Q}^{\hat{\pi}} - \hat{Q}^* + \hat{Q}^{\pi^*} - Q^* \\ &\geq -\left(\|Q^{\hat{\pi}} - \hat{Q}^{\hat{\pi}}\|_\infty + \|\hat{Q}^{\hat{\pi}} - \hat{Q}^*\|_\infty + \|\hat{Q}^{\pi^*} - Q^*\|_\infty\right) \mathbf{1}. \end{aligned} \quad (18)$$

For policy $\hat{\pi}$ satisfying the condition in Theorem 1, we have $\|\hat{Q}^{\hat{\pi}} - \hat{Q}^*\|_\infty \leq \varepsilon_{\text{opt}}$. It boils down to control $\|Q^{\hat{\pi}} - \hat{Q}^{\hat{\pi}}\|_\infty$ and $\|\hat{Q}^{\pi^*} - Q^*\|_\infty$.

To begin with, we can use (11) to further decompose $\|Q^{\hat{\pi}} - \hat{Q}^{\hat{\pi}}\|_{\infty}$ as

$$\begin{aligned}
\|Q^{\hat{\pi}} - \hat{Q}^{\hat{\pi}}\|_{\infty} &= \left\| \left(I - \gamma P^{\hat{\pi}} \right)^{-1} r - \left(I - \gamma \hat{P}^{\hat{\pi}} \right)^{-1} r \right\|_{\infty} \\
&= \left\| \left(I - \gamma P^{\hat{\pi}} \right)^{-1} \left[\left(I - \gamma \hat{P}^{\hat{\pi}} \right) - \left(I - \gamma P^{\hat{\pi}} \right) \right] \hat{Q}^{\hat{\pi}} \right\|_{\infty} \\
&= \left\| \gamma \left(I - \gamma P^{\hat{\pi}} \right)^{-1} (P - \hat{P}) \hat{V}^{\hat{\pi}} \right\|_{\infty} \\
&\leq \left\| \gamma \left(I - \gamma P^{\hat{\pi}} \right)^{-1} (P - \hat{P}) \hat{V}^{\star} \right\|_{\infty} + \left\| \gamma \left(I - \gamma P^{\hat{\pi}} \right)^{-1} (P - \hat{P}) (\hat{V}^{\hat{\pi}} - \hat{V}^{\star}) \right\|_{\infty} \\
&\leq \left\| \gamma \left(I - \gamma P^{\hat{\pi}} \right)^{-1} (P - \hat{P}) \hat{V}^{\star} \right\|_{\infty} + \frac{2\gamma \varepsilon_{\text{opt}}}{1 - \gamma}. \tag{19}
\end{aligned}$$

Here the last inequality is due to

$$\begin{aligned}
&\left\| \gamma \left(I - \gamma P^{\hat{\pi}} \right)^{-1} (P - \hat{P}) (\hat{V}^{\hat{\pi}} - \hat{V}^{\star}) \right\|_{\infty} \\
&\leq \gamma \left\| \left(I - \gamma P^{\hat{\pi}} \right)^{-1} \right\|_1 \left\| (P - \hat{P}) (\hat{V}^{\hat{\pi}} - \hat{V}^{\star}) \right\|_{\infty} \\
&\leq \gamma \left\| \left(I - \gamma P^{\hat{\pi}} \right)^{-1} \right\|_1 (\|P\|_1 + \|\hat{P}\|_1) \|\hat{V}^{\hat{\pi}} - \hat{V}^{\star}\|_{\infty} \\
&\leq \frac{2\gamma \varepsilon_{\text{opt}}}{1 - \gamma},
\end{aligned}$$

where we use the fact that $\|(I - \gamma P^{\hat{\pi}})^{-1}\|_1 \leq 1/(1 - \gamma)$ and $\|P\|_1 = \|\hat{P}\|_1 = 1$.

Similarly, for the term $\|\hat{Q}^{\pi^{\star}} - Q^{\star}\|_{\infty}$ in (18), we have

$$\begin{aligned}
\|\hat{Q}^{\pi^{\star}} - Q^{\star}\|_{\infty} &= \left\| \gamma \left(I - \gamma P^{\pi^{\star}} \right)^{-1} (P - \hat{P}) \hat{V}^{\pi^{\star}} \right\|_{\infty} \\
&\leq \left\| \gamma \left(I - \gamma P^{\pi^{\star}} \right)^{-1} (P - \hat{P}) \hat{V}^{\pi^{\star}} \right\|_{\infty}. \tag{20}
\end{aligned}$$

As can be seen from (19) and (20), it boils down to bound $|(P - \hat{P})\hat{V}^{\star}|$ and $|(P - \hat{P})\hat{V}^{\pi^{\star}}|$. We have the following lemma.

Lemma 1. *With probability exceeding $1 - \delta$, one has*

$$\begin{aligned}
\left| (P - \hat{P})_{s,a} \hat{V}^{\star} \right| &\leq \frac{10\xi}{1 - \gamma} + 4\sqrt{\frac{2\log(4K/\delta)}{N}} + \frac{4\log(8K/((1 - \gamma)\delta))}{(1 - \gamma)N} \\
&\quad + \sqrt{\frac{4\log(8K/((1 - \gamma)\delta))}{N}} \sqrt{\text{Var}_{P_{s,a}}(\hat{V}^{\star})}, \tag{21}
\end{aligned}$$

$$\begin{aligned}
\left| (P - \hat{P})_{s,a} \hat{V}^{\pi^{\star}} \right| &\leq \frac{10\xi}{1 - \gamma} + 4\sqrt{\frac{2\log(4K/\delta)}{N}} + \frac{4\log(8K/((1 - \gamma)\delta))}{(1 - \gamma)N} \\
&\quad + \sqrt{\frac{4\log(8K/((1 - \gamma)\delta))}{N}} \sqrt{\text{Var}_{P_{s,a}}(\hat{V}^{\pi^{\star}})}. \tag{22}
\end{aligned}$$

Proof. See Appendix B.2. □

Applying (21) to (19) reveals that

$$\begin{aligned} \left\| \mathbf{Q}^{\hat{\pi}} - \hat{\mathbf{Q}}^{\hat{\pi}} \right\|_{\infty} &\leq \sqrt{\frac{4 \log(8K/((1-\gamma)\delta))}{N}} \left\| \gamma \left(\mathbf{I} - \gamma \mathbf{P}^{\hat{\pi}} \right)^{-1} \sqrt{\text{Var}_{\mathbf{P}_{s,a}}(\hat{\mathbf{V}}^{\star})} \right\|_{\infty} \\ &\quad + \frac{\gamma}{1-\gamma} \left[4 \sqrt{\frac{2 \log(4K/\delta)}{N}} + \frac{4 \log(8K/((1-\gamma)\delta))}{(1-\gamma)N} \right] \\ &\quad + \frac{10\gamma\xi}{(1-\gamma)^2} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}. \end{aligned} \quad (23)$$

For the first term, one has

$$\begin{aligned} \sqrt{\text{Var}_{\mathbf{P}_{s,a}}(\hat{\mathbf{V}}^{\star})} &\leq \sqrt{\text{Var}_{\mathbf{P}_{s,a}}(\mathbf{V}^{\hat{\pi}})} + \sqrt{\text{Var}_{\mathbf{P}_{s,a}}(\mathbf{V}^{\hat{\pi}} - \hat{\mathbf{V}}^{\hat{\pi}})} + \sqrt{\text{Var}_{\mathbf{P}_{s,a}}(\hat{\mathbf{V}}^{\hat{\pi}} - \hat{\mathbf{V}}^{\star})} \\ &\leq \sqrt{\text{Var}_{\mathbf{P}_{s,a}}(\mathbf{V}^{\hat{\pi}})} + \left\| \mathbf{V}^{\hat{\pi}} - \hat{\mathbf{V}}^{\hat{\pi}} \right\|_{\infty} + \varepsilon_{\text{opt}} \\ &\leq \sqrt{\text{Var}_{\mathbf{P}_{s,a}}(\mathbf{V}^{\hat{\pi}})} + \left\| \mathbf{Q}^{\hat{\pi}} - \hat{\mathbf{Q}}^{\hat{\pi}} \right\|_{\infty} + \varepsilon_{\text{opt}}, \end{aligned}$$

where the first inequality comes from the fact that $\sqrt{\text{Var}(X+Y)} \leq \sqrt{\text{Var}(X)} + \sqrt{\text{Var}(Y)}$ for any random variables X and Y . It follows that

$$\begin{aligned} &\left\| \gamma \left(\mathbf{I} - \gamma \mathbf{P}^{\hat{\pi}} \right)^{-1} \sqrt{\text{Var}_{\mathbf{P}_{s,a}}(\hat{\mathbf{V}}^{\star})} \right\|_{\infty} \\ &\leq \left\| \gamma \left(\mathbf{I} - \gamma \mathbf{P}^{\hat{\pi}} \right)^{-1} \sqrt{\text{Var}_{\mathbf{P}_{s,a}}(\mathbf{V}^{\hat{\pi}})} \right\|_{\infty} + \frac{\gamma}{1-\gamma} \left(\left\| \mathbf{Q}^{\hat{\pi}} - \hat{\mathbf{Q}}^{\hat{\pi}} \right\|_{\infty} + \varepsilon_{\text{opt}} \right) \\ &\leq \gamma \sqrt{\frac{2}{(1-\gamma)^3}} + \frac{\gamma}{1-\gamma} \left(\left\| \mathbf{Q}^{\hat{\pi}} - \hat{\mathbf{Q}}^{\hat{\pi}} \right\|_{\infty} + \varepsilon_{\text{opt}} \right), \end{aligned} \quad (24)$$

where the second inequality utilizes [3, Lemma 7].

Plugging (24) into (23) yields

$$\begin{aligned} \left\| \mathbf{Q}^{\hat{\pi}} - \hat{\mathbf{Q}}^{\hat{\pi}} \right\|_{\infty} &\leq \sqrt{\frac{4 \log(8K/((1-\gamma)\delta))}{N}} \left[\gamma \sqrt{\frac{2}{(1-\gamma)^3}} + \frac{\gamma}{1-\gamma} \left(\left\| \mathbf{Q}^{\hat{\pi}} - \hat{\mathbf{Q}}^{\hat{\pi}} \right\|_{\infty} + \varepsilon_{\text{opt}} \right) \right] \\ &\quad + \frac{\gamma}{1-\gamma} \left[4 \sqrt{\frac{2 \log(4K/\delta)}{N}} + \frac{4 \log(8K/((1-\gamma)\delta))}{(1-\gamma)N} \right] + \frac{10\gamma\xi}{(1-\gamma)^2} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}. \end{aligned}$$

Then we can rearrange terms to obtain

$$\left\| \mathbf{Q}^{\hat{\pi}} - \hat{\mathbf{Q}}^{\hat{\pi}} \right\|_{\infty} \leq 10\gamma \sqrt{\frac{\log(8K/((1-\gamma)\delta))}{N(1-\gamma)^3}} + \frac{11\gamma\xi}{(1-\gamma)^2} + \frac{3\gamma\varepsilon_{\text{opt}}}{1-\gamma} \quad (25)$$

as long as $N \geq \tilde{C} \log(8K/((1-\gamma)\delta))/(1-\gamma)^2$ for some sufficiently large constant $\tilde{C} > 0$.

In a similar vein, we can use (20) and (22) to obtain that

$$\left\| \hat{\mathbf{Q}}^{\pi^{\star}} - \mathbf{Q}^{\star} \right\|_{\infty} \leq 10\gamma \sqrt{\frac{\log(8K/((1-\gamma)\delta))}{N(1-\gamma)^3}} + \frac{11\gamma\xi}{(1-\gamma)^2}. \quad (26)$$

Finally, we can substitute (25) and (26) into (18) to achieve

$$\mathbf{Q}^{\hat{\pi}} - \mathbf{Q}^{\star} \geq - \left(20\gamma \sqrt{\frac{\log(8K/((1-\gamma)\delta))}{N(1-\gamma)^3}} + \frac{22\gamma\xi}{(1-\gamma)^2} + \frac{3\gamma\varepsilon_{\text{opt}}}{1-\gamma} + \varepsilon_{\text{opt}} \right) \mathbf{1}.$$

This result implies that

$$\mathbf{Q}^{\hat{\pi}} \geq \mathbf{Q}^{\star} - \left(\varepsilon + \frac{22\xi}{(1-\gamma)^2} + \frac{4\varepsilon_{\text{opt}}}{1-\gamma} \right) \mathbf{1},$$

as long as

$$N \geq \frac{C \log(8K/((1-\gamma)\delta))}{(1-\gamma)^3 \varepsilon^2},$$

for some sufficiently large constant $C > 0$.

B.2 Proof of Lemma 1

To prove this theorem, we invoke the idea of s -absorbing MDP proposed by [1]. For a state $s \in \mathcal{S}$ and a scalar u , we define a new MDP $M_{s,u}$ to be identical to M on all the other states except s ; on state s , $M_{s,u}$ is absorbing such that $P_{M_{s,u}}(s|s, a) = 1$ and $r_{M_{s,u}}(s, a) = (1 - \gamma)u$ for all $a \in \mathcal{A}$. More formally, we define $P_{M_{s,u}}$ and $r_{M_{s,u}}$ as

$$\begin{aligned} P_{M_{s,u}}(s|s, a) &= 1, \quad r_{M_{s,u}}(s, a) = (1 - \gamma)u, \quad \text{for all } a \in \mathcal{A}, \\ P_{M_{s,u}}(\cdot|s', a') &= P(\cdot|s', a'), \quad r_{M_{s,u}}(s, a) = r(s, a), \quad \text{for all } s' \neq s \text{ and } a' \in \mathcal{A}. \end{aligned}$$

To streamline notations, we will use $\mathbf{V}_{s,u}^\pi \in \mathbb{R}^{|\mathcal{S}|}$ and $\mathbf{V}_{s,u}^* \in \mathbb{R}^{|\mathcal{S}|}$ to denote the value function of $M_{s,u}$ under policy π and the optimal value function of $M_{s,u}$ respectively. Furthermore, we denote by $\widehat{M}_{s,u}$ the MDP whose probability transition kernel is identical to \widehat{P} at all states except that state s is absorbing. Similar as before, we use $\widehat{\mathbf{V}}_{s,u}^* \in \mathbb{R}^{|\mathcal{S}|}$ to denote the optimal value function under $\widehat{M}_{s,u}$. The construction of this collection of auxiliary MDPs will facilitate our analysis by decoupling the statistical dependency between \widehat{P} and $\widehat{\pi}^*$.

To begin with, we can decompose the quantity of interest as

$$\begin{aligned} \left| \left(\mathbf{P} - \widehat{\mathbf{P}} \right)_{s,a} \widehat{\mathbf{V}}^* \right| &= \left| \left(\mathbf{P} - \widehat{\mathbf{P}} \right)_{s,a} \left(\widehat{\mathbf{V}}^* - \widehat{\mathbf{V}}_{s,u}^* + \widehat{\mathbf{V}}_{s,u}^* \right) \right| \\ &\leq \left| \left(\mathbf{P} - \widehat{\mathbf{P}} \right)_{s,a} \widehat{\mathbf{V}}_{s,u}^* \right| + \left| \left(\mathbf{P} - \widehat{\mathbf{P}} \right)_{s,a} \left(\widehat{\mathbf{V}}^* - \widehat{\mathbf{V}}_{s,u}^* \right) \right| \\ &\stackrel{(i)}{\leq} \left| \left(\mathbf{P} - \widetilde{\mathbf{P}} \right)_{s,a} \widehat{\mathbf{V}}_{s,u}^* \right| + \left| \lambda(s, a) \left(\widetilde{\mathbf{P}}_{\mathcal{K}} - \mathbf{P}_{\mathcal{K}} \right) \widehat{\mathbf{V}}_{s,u}^* \right| \\ &\quad + \left| \lambda(s, a) \left(\mathbf{P}_{\mathcal{K}} - \widehat{\mathbf{P}}_{\mathcal{K}} \right) \widehat{\mathbf{V}}_{s,u}^* \right| + \left(\|\mathbf{P}_{s,a}\|_1 + \|\widehat{\mathbf{P}}_{s,a}\|_1 \right) \|\widehat{\mathbf{V}}^* - \widehat{\mathbf{V}}_{s,u}^*\|_\infty \\ &\leq \left\| \left(\mathbf{P} - \widetilde{\mathbf{P}} \right)_{s,a} \right\|_1 \|\widehat{\mathbf{V}}_{s,u}^*\|_\infty + \|\lambda(s, a)\|_1 \cdot \left\| \left(\widetilde{\mathbf{P}}_{\mathcal{K}} - \mathbf{P}_{\mathcal{K}} \right) \widehat{\mathbf{V}}_{s,u}^* \right\|_\infty \\ &\quad + \|\lambda(s, a)\|_1 \cdot \left\| \left(\mathbf{P}_{\mathcal{K}} - \widehat{\mathbf{P}}_{\mathcal{K}} \right) \widehat{\mathbf{V}}_{s,u}^* \right\|_\infty + 2 \|\widehat{\mathbf{V}}^* - \widehat{\mathbf{V}}_{s,u}^*\|_\infty \\ &\stackrel{(ii)}{\leq} \frac{2\xi}{1 - \gamma} + \max_{(s,a) \in \mathcal{K}} \left| \left(\mathbf{P} - \widehat{\mathbf{P}} \right)_{s,a} \widehat{\mathbf{V}}_{s,u}^* \right| + 2 \|\widehat{\mathbf{V}}^* - \widehat{\mathbf{V}}_{s,u}^*\|_\infty, \end{aligned} \quad (27)$$

where (i) makes use of $\widetilde{\mathbf{P}}_{s,a} = \lambda(s, a) \widetilde{\mathbf{P}}_{\mathcal{K}}$ and $\widehat{\mathbf{P}}_{s,a} = \lambda(s, a) \widehat{\mathbf{P}}_{\mathcal{K}}$; (ii) depends on $\|\mathbf{P} - \widetilde{\mathbf{P}}\|_1 \leq \xi$, $\|\lambda(s, a)\|_1 = 1$ and $\|\widehat{\mathbf{V}}_{s,u}^*\|_\infty \leq (1 - \gamma)^{-1}$. For each state s , the value of u will be selected from a set \mathcal{U}_s . The choice of \mathcal{U}_s will be specified later. Then for some fixed u in \mathcal{U}_s and fixed state-action pair $(s, a) \in \mathcal{K}$, due to the independence between $\widehat{\mathbf{P}}_{s,a}$ and $\widehat{\mathbf{V}}_{s,u}^*$, we can apply Bernstein's inequality (cf. [5, Theorem 2.8.4]) conditional on $\widehat{\mathbf{V}}_{s,u}^*$ to reveal that with probability greater than $1 - \delta/2$,

$$\left| \left(\mathbf{P} - \widehat{\mathbf{P}} \right)_{s,a} \widehat{\mathbf{V}}_{s,u}^* \right| \leq \sqrt{\frac{2 \log(4/\delta)}{N} \text{Var}_{\mathbf{P}_{s,a}} \left(\widehat{\mathbf{V}}_{s,u}^* \right)} + \frac{2 \log(4/\delta)}{3(1 - \gamma)N}. \quad (28)$$

Invoking the union bound over all the K state-action pairs of \mathcal{K} and all the possible values of u in \mathcal{U}_s demonstrate that with probability greater than $1 - \delta/2$,

$$\left| \left(\mathbf{P} - \widehat{\mathbf{P}} \right)_{s,a} \widehat{\mathbf{V}}_{s,u}^* \right| \leq \sqrt{\frac{2 \log(4K|\mathcal{U}_s|/\delta)}{N} \text{Var}_{\mathbf{P}_{s,a}} \left(\widehat{\mathbf{V}}_{s,u}^* \right)} + \frac{2 \log(4K|\mathcal{U}_s|/\delta)}{3(1 - \gamma)N}, \quad (29)$$

holds for all state-action pair $(s, a) \in \mathcal{K}$ and all $u \in \mathcal{U}_s$. Here, $\text{Var}_{\mathbf{P}_{s,a}}(\cdot)$ is defined in (13). Then we observe that

$$\begin{aligned} \sqrt{\text{Var}_{\mathbf{P}_{s,a}} \left(\widehat{\mathbf{V}}_{s,u}^* \right)} &\leq \sqrt{\text{Var}_{\mathbf{P}_{s,a}} \left(\widehat{\mathbf{V}}^* - \widehat{\mathbf{V}}_{s,u}^* \right)} + \sqrt{\text{Var}_{\mathbf{P}_{s,a}} \left(\widehat{\mathbf{V}}^* \right)} \\ &\leq \|\widehat{\mathbf{V}}^* - \widehat{\mathbf{V}}_{s,u}^*\|_\infty + \sqrt{\text{Var}_{\mathbf{P}_{s,a}} \left(\widehat{\mathbf{V}}^* \right)} \\ &\leq \left| \widehat{\mathbf{V}}^*(s) - u \right| + \sqrt{\text{Var}_{\mathbf{P}_{s,a}} \left(\widehat{\mathbf{V}}^* \right)}, \end{aligned} \quad (30)$$

where (i) is due to $\sqrt{\text{Var}_{\mathbf{P}_{s,a}}(\mathbf{V}_1 + \mathbf{V}_2)} \leq \sqrt{\text{Var}_{\mathbf{P}_{s,a}}(\mathbf{V}_1)} + \sqrt{\text{Var}_{\mathbf{P}_{s,a}}(\mathbf{V}_2)}$ and (ii) holds since

$$\left\| \hat{\mathbf{V}}^* - \hat{\mathbf{V}}_{s,u}^* \right\|_{\infty} = \left\| \hat{\mathbf{V}}_{s,\hat{\mathbf{V}}^*(s)}^* - \hat{\mathbf{V}}_{s,u}^* \right\|_{\infty} \leq \left| \hat{\mathbf{V}}^*(s) - u \right|, \quad (31)$$

whose proof can be found in [1, Lemma 8 and 9].

By substituting (29), (30) and (31) into (27), we arrive at

$$\begin{aligned} \left| (\mathbf{P} - \hat{\mathbf{P}})_{s,a} \hat{\mathbf{V}}^* \right| &\leq \frac{2\xi}{1-\gamma} + \left| \hat{\mathbf{V}}^*(s) - u \right| \left(2 + \sqrt{\frac{2 \log(4K|\mathcal{U}_s|/\delta)}{N}} \right) \\ &\quad + \sqrt{\frac{2 \log(4K|\mathcal{U}_s|/\delta)}{N}} \sqrt{\text{Var}_{\mathbf{P}_{s,a}}(\hat{\mathbf{V}}^*)} + \frac{2 \log(4K|\mathcal{U}_s|/\delta)}{3(1-\gamma)N}. \end{aligned} \quad (32)$$

Then it boils down to determining \mathcal{U}_s . The coarse bounds of $\hat{\mathbf{Q}}^{\pi^*}$ and $\hat{\mathbf{Q}}^*$ in the following lemma provide a guidance on the choice of \mathcal{U}_s .

Lemma 2. For $\delta \in (0, 1)$, with probability exceeding $1 - \delta/2$ one has

$$\left\| \mathbf{Q}^* - \hat{\mathbf{Q}}^{\pi^*} \right\|_{\infty} \leq \frac{\gamma}{1-\gamma} \sqrt{\frac{\log(4K/\delta)}{2N(1-\gamma)^2}} + \frac{2\gamma\xi}{(1-\gamma)^2}, \quad (33)$$

$$\left\| \mathbf{Q}^* - \hat{\mathbf{Q}}^* \right\|_{\infty} \leq \frac{\gamma}{1-\gamma} \sqrt{\frac{\log(4K/\delta)}{2N(1-\gamma)^2}} + \frac{2\gamma\xi}{(1-\gamma)^2}. \quad (34)$$

Proof. See Appendix B.3. □

This inspires us to choose \mathcal{U}_s to be the set consisting of equidistant points in $[\mathbf{V}^*(s) - R(\delta), \mathbf{V}^*(s) + R(\delta)]$ with $|\mathcal{U}_s| = \lceil 1/(1-\gamma)^2 \rceil$ and

$$R(\delta) := \frac{\gamma}{1-\gamma} \sqrt{\frac{\log(4K/\delta)}{2N(1-\gamma)^2}} + \frac{2\gamma\xi}{(1-\gamma)^2}.$$

Since $\|\mathbf{V}^* - \hat{\mathbf{V}}^*\|_{\infty} \leq \|\mathbf{Q}^* - \hat{\mathbf{Q}}^*\|_{\infty}$, Lemma 2 implies that $\hat{\mathbf{V}}^*(s) \in [\mathbf{V}^*(s) - R(\delta), \mathbf{V}^*(s) + R(\delta)]$ with probability over $1 - \delta/2$. Hence, we have

$$\min_{u \in \mathcal{U}_s} \left| \hat{\mathbf{V}}^*(s) - u \right| \leq \frac{2R(\delta)}{|\mathcal{U}_s| + 1} \leq 2\gamma \sqrt{\frac{2 \log(4K/\delta)}{N}} + 4\gamma\xi. \quad (35)$$

Consequently, with probability exceeding $1 - \delta$, one has

$$\begin{aligned} \left| (\mathbf{P} - \hat{\mathbf{P}})_{s,a} \hat{\mathbf{V}}^* \right| &\stackrel{(i)}{\leq} \frac{2\xi}{1-\gamma} + \min_{u \in \mathcal{U}_s} \left| \hat{\mathbf{V}}^*(s) - u \right| \left(2 + \sqrt{\frac{2 \log(4K|\mathcal{U}_s|/\delta)}{N}} \right) \\ &\quad + \sqrt{\frac{2 \log(4K|\mathcal{U}_s|/\delta)}{N}} \sqrt{\text{Var}_{\mathbf{P}_{s,a}}(\hat{\mathbf{V}}^*)} + \frac{2 \log(4K|\mathcal{U}_s|/\delta)}{3(1-\gamma)N} \\ &\stackrel{(ii)}{\leq} \frac{2\xi}{1-\gamma} + \left(2\gamma \sqrt{\frac{2 \log(4K/\delta)}{N}} + 4\gamma\xi \right) \left(2 + \sqrt{\frac{4 \log(8K/((1-\gamma)\delta))}{N}} \right) \\ &\quad + \sqrt{\frac{4 \log(8K/((1-\gamma)\delta))}{N}} \sqrt{\text{Var}_{\mathbf{P}_{s,a}}(\hat{\mathbf{V}}^*)} + \frac{2 \log(8K/((1-\gamma)\delta))}{3(1-\gamma)N} \\ &\leq \frac{10\xi}{1-\gamma} + 4\sqrt{\frac{2 \log(4K/\delta)}{N}} + \frac{4 \log(8K/((1-\gamma)\delta))}{(1-\gamma)N} \\ &\quad + \sqrt{\frac{4 \log(8K/((1-\gamma)\delta))}{N}} \sqrt{\text{Var}_{\mathbf{P}_{s,a}}(\hat{\mathbf{V}}^*)}, \end{aligned}$$

where (i) follows from (32) and (ii) utilizes (35). This finishes the proof for the first inequality. The second inequality can be proved in a similar way and is omitted here for brevity.

B.3 Proof of Lemma 2

To begin with, one has

$$\begin{aligned}
\|(\hat{P} - P) V^*\|_\infty &\leq \|\Lambda(\hat{P}_K - P_K) V^*\|_\infty + \|\Lambda(P_K - \tilde{P}_K) V^*\|_\infty + \|(\tilde{P} - P) V^*\|_\infty \\
&\leq \|\Lambda\|_1 \|(\hat{P}_K - P_K) V^*\|_\infty + \|\Lambda\|_1 \|(P_K - \tilde{P}_K) V^*\|_\infty + \|\tilde{P} - P\|_1 \|V^*\|_\infty \\
&\leq \|(\hat{P}_K - P_K) V^*\|_\infty + \frac{2\xi}{1-\gamma},
\end{aligned} \tag{36}$$

where the first line uses $\hat{P} = \Lambda \hat{P}_K$ and $\tilde{P} = \Lambda \tilde{P}_K$; the last inequality comes from the facts that $\|\tilde{P} - P\|_1 \leq \xi$, $\|\Lambda\|_1 = 1$ and $\|V^*\|_\infty \leq (1-\gamma)^{-1}$. Then we turn to bound $\|(\hat{P}_K - P_K) V^*\|_\infty$. In view of (4), Hoeffding's inequality (cf. [5, Theorem 2.2.6]) implies that for $(s, a) \in \mathcal{K}$,

$$\mathbb{P}\left(\left|(\hat{P} - P)_{s,a} V^*\right| \geq t\right) \leq 2 \exp\left(-\frac{2t^2}{\|V^*\|_\infty^2 / N}\right).$$

Hence by the standard union bound argument we have

$$\|(\hat{P}_K - P_K) V^*\|_\infty \leq \sqrt{\frac{\|V^*\|_\infty^2 \log(4K/\delta)}{2N}} \leq \sqrt{\frac{\log(4K/\delta)}{2N(1-\gamma)^2}}, \tag{37}$$

with probability over $1 - \delta/2$.

1. Now we are ready to bound $Q^{\pi^*} - \hat{Q}^{\pi^*}$. One has

$$\begin{aligned}
Q^{\pi^*} - \hat{Q}^{\pi^*} &= (I - \gamma P^{\pi^*})^{-1} r - (I - \gamma \hat{P}^{\pi^*})^{-1} r \\
&= (I - \gamma \hat{P}^{\pi^*})^{-1} \left((I - \gamma \hat{P}^{\pi^*}) - (I - \gamma P^{\pi^*}) \right) Q^{\pi^*} \\
&= \gamma (I - \gamma \hat{P}^{\pi^*})^{-1} (P^{\pi^*} - \hat{P}^{\pi^*}) Q^{\pi^*} \\
&= \gamma (I - \gamma \hat{P}^{\pi^*})^{-1} (P - \hat{P}) V^{\pi^*},
\end{aligned}$$

where the first equality makes use of (11). Then we take (36) and (37) collectively to achieve

$$\begin{aligned}
\left\| \gamma (I - \gamma \hat{P}^{\pi^*})^{-1} (P - \hat{P}) V^* \right\|_\infty &\leq \gamma \sum_{i=0}^{\infty} \left\| \gamma^i (\hat{P}^{\pi^*})^i (P - \hat{P}) V^* \right\|_\infty \\
&\leq \gamma \sum_{i=0}^{\infty} \gamma^i \left\| (\hat{P}^{\pi^*})^i \right\|_1 \left\| (P - \hat{P}) V^* \right\|_\infty \\
&\leq \frac{\gamma}{1-\gamma} \sqrt{\frac{\log(4K/\delta)}{2N(1-\gamma)^2}} + \frac{2\gamma\xi}{(1-\gamma)^2},
\end{aligned}$$

where the last line comes from the fact that for all $i \geq 1$, $(\hat{P}^{\pi^*})^i$ is a probability transition matrix so that $\|(\hat{P}^{\pi^*})^i\|_1 = 1$. This justifies the first inequality (33).

2. In terms of the second one, [1, Section A.4] implies that

$$\left\| Q^* - \hat{Q}^* \right\|_\infty \leq \frac{\gamma}{1-\gamma} \left\| (P - \hat{P}) V^* \right\|_\infty.$$

Substitution of (36) and (37) into the above inequality yields

$$\left\| Q^* - \hat{Q}^* \right\|_\infty \leq \frac{\gamma}{1-\gamma} \sqrt{\frac{\log(4K/\delta)}{2N(1-\gamma)^2}} + \frac{2\gamma\xi}{(1-\gamma)^2}.$$

C Analysis of Q-learning (Proof of Theorem 2)

In this section, we will provide complete proof for Theorem 2. We actually prove a more general version of Theorem 2 that takes model misspecification into consideration, as stated below.

Theorem 4. Consider any $\delta \in (0, 1)$ and $\varepsilon \in (0, 1]$. Suppose that there exists a probability transition model \tilde{P} obeying Definition 1 and Assumption 1 with feature vectors $\{\phi(s, a)\}_{(s,a) \in \mathcal{S} \times \mathcal{A}} \subset \mathbb{R}^K$ and anchor state-action pairs \mathcal{K} such that

$$\|\tilde{P} - P\|_1 \leq \xi$$

for some $\xi \geq 0$. Assume that the initialization obeys $0 \leq Q_0(s, a) \leq \frac{1}{1-\gamma}$ for any $(s, a) \in \mathcal{S} \times \mathcal{A}$ and for any $0 \leq t \leq T$, the learning rates satisfy

$$\frac{1}{1 + \frac{c_1(1-\gamma)T}{\log^2 T}} \leq \eta_t \leq \frac{1}{1 + \frac{c_2(1-\gamma)t}{\log^2 T}}, \quad (38)$$

for some sufficiently small universal constants $c_1 \geq c_2 > 0$. Suppose that the total number of iterations T exceeds

$$T \geq \frac{C_3 \log(KT/\delta) \log^4 T}{(1-\gamma)^4 \varepsilon^2}, \quad (39)$$

for some sufficiently large universal constant $C_3 > 0$. If there exists a linear probability transition model \tilde{P} satisfying Assumption 1 with feature vectors $\{\phi(s, a)\}_{(s,a) \in \mathcal{S} \times \mathcal{A}}$ such that $\|\tilde{P} - P\|_1 \leq \xi$, then with probability exceeding $1 - \delta$, the output Q_T of Algorithm 2 satisfies

$$\max_{(s,a) \in \mathcal{S} \times \mathcal{A}} |Q_T(s, a) - Q^*(s, a)| \leq \varepsilon + \frac{6\gamma\xi}{(1-\gamma)^2}, \quad (40)$$

for some constant $C_4 > 0$. In addition, let π_T (resp. V_T) to be the policy (resp. value function) induced by Q_T , then one has

$$\max_{s \in \mathcal{S}} |V^{\pi_T}(s) - V^*(s)| \leq \frac{2\gamma}{1-\gamma} \left(\varepsilon + \frac{6\gamma\xi}{(1-\gamma)^2} \right). \quad (41)$$

Theorem 4 subsumes Theorem 2 as a special case with $\xi = 0$. The remainder of this section is devoted to proving Theorem 4.

C.1 Proof of Theorem 4

First we show that (41) can be easily obtained from (40). Since [49] gives rise to

$$\|V^{\pi_T} - V^*\|_\infty \leq \frac{2\gamma\|V_T - V^*\|_\infty}{1-\gamma},$$

we have

$$\|V^{\pi_T} - V^*\|_\infty \leq \frac{2\gamma\|Q_T - Q^*\|_\infty}{1-\gamma},$$

due to $\|V_T - V^*\|_\infty \leq \|Q_T - Q^*\|_\infty$. Then (41) follows directly from (40).

Therefore, we are left to justify (40). To start with, we consider the update rule

$$Q_t = (1 - \eta_t) Q_{t-1} + \eta_t (r + \gamma \hat{P}_t V_{t-1}).$$

By defining the error term $\Delta_t := Q_t - Q^*$, we can decompose Δ_t into

$$\begin{aligned} \Delta_t &= (1 - \eta_t) Q_{t-1} + \eta_t (r + \gamma \hat{P}_t V_{t-1}) - Q^* \\ &= (1 - \eta_t) (Q_{t-1} - Q^*) + \eta_t (r + \gamma \hat{P}_t V_{t-1} - Q^*) \\ &= (1 - \eta_t) (Q_{t-1} - Q^*) + \gamma \eta_t (\hat{P}_t V_{t-1} - P V^*) \\ &= (1 - \eta_t) \Delta_{t-1} + \gamma \eta_t \Lambda (\hat{P}_K^{(t)} - P_K) V_{t-1} + \gamma \eta_t \Lambda P_K (V_{t-1} - V^*) \\ &\quad + \gamma \eta_t (\Lambda P_K - P) V^*. \end{aligned} \quad (42)$$

Here in the penultimate equality, we make use of $\mathbf{Q}^* = \mathbf{r} + \gamma \mathbf{P} \mathbf{V}^*$; and the last equality comes from $\hat{\mathbf{P}}_t = \mathbf{\Lambda} \hat{\mathbf{P}}_{\mathcal{K}}^{(t)}$ which is defined in (15). It is straightforward to check that $\mathbf{\Lambda} \mathbf{P}_{\mathcal{K}}$ is also a probability transition matrix. We denote by $\bar{\mathbf{P}} = \mathbf{\Lambda} \mathbf{P}_{\mathcal{K}}$ hereafter. The third term in the decomposition above can be upper and lower bounded by

$$\bar{\mathbf{P}}(\mathbf{V}_{t-1} - \mathbf{V}^*) = \bar{\mathbf{P}}^{\pi_{t-1}} \mathbf{Q}_{t-1} - \bar{\mathbf{P}}^{\pi^*} \mathbf{Q}^* \leq \bar{\mathbf{P}}^{\pi_{t-1}} \mathbf{Q}_{t-1} - \bar{\mathbf{P}}^{\pi_{t-1}} \mathbf{Q}^* = \bar{\mathbf{P}}^{\pi_{t-1}} \Delta_{t-1},$$

and

$$\bar{\mathbf{P}}(\mathbf{V}_{t-1} - \mathbf{V}^*) = \bar{\mathbf{P}}^{\pi_{t-1}} \mathbf{Q}_{t-1} - \bar{\mathbf{P}}^{\pi^*} \mathbf{Q}^* \geq \bar{\mathbf{P}}^{\pi^*} \mathbf{Q}_{t-1} - \bar{\mathbf{P}}^{\pi^*} \mathbf{Q}^* = \bar{\mathbf{P}}^{\pi^*} \Delta_{t-1}.$$

Plugging these bounds into (42) yields

$$\Delta_t \leq (1 - \eta_t) \Delta_{t-1} + \gamma \eta_t \mathbf{\Lambda} \left(\hat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{t-1} + \gamma \eta_t \bar{\mathbf{P}}^{\pi_{t-1}} \Delta_{t-1} + \gamma \eta_t (\mathbf{\Lambda} \mathbf{P}_{\mathcal{K}} - \mathbf{P}) \mathbf{V}^*,$$

$$\Delta_t \geq (1 - \eta_t) \Delta_{t-1} + \gamma \eta_t \mathbf{\Lambda} \left(\hat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{t-1} + \gamma \eta_t \bar{\mathbf{P}}^{\pi^*} \Delta_{t-1} + \gamma \eta_t (\mathbf{\Lambda} \mathbf{P}_{\mathcal{K}} - \mathbf{P}) \mathbf{V}^*.$$

Repeatedly invoking these two recursive relations leads to

$$\Delta_t \leq \eta_0^{(t)} \Delta_0 + \sum_{i=1}^t \eta_i^{(t)} \gamma \left(\bar{\mathbf{P}}^{\pi_{t-1}} \Delta_{t-1} + \mathbf{\Lambda} \left(\hat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{t-1} + (\mathbf{\Lambda} \mathbf{P}_{\mathcal{K}} - \mathbf{P}) \mathbf{V}^* \right), \quad (43)$$

$$\Delta_t \geq \eta_0^{(t)} \Delta_0 + \sum_{i=1}^t \eta_i^{(t)} \gamma \left(\bar{\mathbf{P}}^{\pi^*} \Delta_{t-1} + \mathbf{\Lambda} \left(\hat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{t-1} + (\mathbf{\Lambda} \mathbf{P}_{\mathcal{K}} - \mathbf{P}) \mathbf{V}^* \right), \quad (44)$$

where

$$\eta_i^{(t)} := \begin{cases} \prod_{j=1}^t (1 - \eta_j), & \text{if } i = 0, \\ \eta_i \prod_{j=i+1}^t (1 - \eta_j), & \text{if } 0 < i < t, \\ \eta_t, & \text{if } i = t. \end{cases}$$

Here we adopt the same notations as [4].

To begin with, we consider the upper bound (43). It can be further decomposed as

$$\begin{aligned} \Delta_t &\leq \underbrace{\eta_0^{(t)} \Delta_0 + \sum_{i=1}^{(1-\alpha)t} \eta_i^{(t)} \gamma \left(\bar{\mathbf{P}}^{\pi_{t-1}} \Delta_{t-1} + \mathbf{\Lambda} \left(\hat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{t-1} \right)}_{=:\boldsymbol{\theta}_t} \\ &\quad + \underbrace{\sum_{i=(1-\alpha)t+1}^t \eta_i^{(t)} \gamma \mathbf{\Lambda} \left(\hat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{i-1}}_{=:\boldsymbol{\nu}_t} \\ &\quad + \underbrace{\sum_{i=1}^t \eta_i^{(t)} \gamma (\mathbf{\Lambda} \mathbf{P}_{\mathcal{K}} - \mathbf{P}) \mathbf{V}^*}_{=:\boldsymbol{\omega}_t} + \sum_{i=(1-\alpha)t+1}^t \eta_i^{(t)} \gamma \bar{\mathbf{P}}^{\pi_{t-1}} \Delta_{i-1}, \end{aligned} \quad (45)$$

where we define $\alpha := C_4(1 - \gamma)/\log T$ for some constant $C_4 > 0$. Next, we turn to bound $\boldsymbol{\theta}_t$ and $\boldsymbol{\nu}_t$ respectively for any t satisfying $\frac{T}{c_2 \log \frac{1}{1-\gamma}} \leq t \leq T$ with stepsize choice (8).

Bounding $\boldsymbol{\omega}_t$. It is straightforward to bound

$$\begin{aligned} \|\boldsymbol{\omega}_t\|_{\infty} &\stackrel{(i)}{=} \|\gamma (\mathbf{\Lambda} \mathbf{P}_{\mathcal{K}} - \mathbf{P}) \mathbf{V}^*\|_{\infty} \\ &\stackrel{(ii)}{\leq} \gamma \left(\|\mathbf{\Lambda}\|_1 \left\| (\mathbf{P}_{\mathcal{K}} - \tilde{\mathbf{P}}_{\mathcal{K}}) \mathbf{V}^* \right\|_{\infty} + \left\| (\tilde{\mathbf{P}} - \mathbf{P}) \mathbf{V}^* \right\|_{\infty} \right) \\ &\stackrel{(iii)}{\leq} \frac{2\gamma\xi}{1-\gamma}, \end{aligned}$$

where the first equality comes from the fact that $\sum_{i=1}^t \eta_i^{(t)} = 1$ [4, Equation (40)]; the second inequality utilizes $\tilde{\mathbf{P}} = \mathbf{\Lambda} \tilde{\mathbf{P}}_{\mathcal{K}}$; the last line uses the facts that $\|\mathbf{\Lambda}\|_1 = 1$, $\|\mathbf{V}^*\|_{\infty} \leq (1 - \gamma)^{-1}$ and $\|\tilde{\mathbf{P}}_{\mathcal{K}} - \mathbf{P}_{\mathcal{K}}\|_1 \leq \|\tilde{\mathbf{P}} - \mathbf{P}\|_1 \leq \xi$.

Bounding θ_t . By similar derivation as Step 1 in [4, Appendix A.2], we have

$$\begin{aligned}
\|\theta_t\|_\infty &\leq \eta_0^{(t)} \|\Delta_0\|_\infty + t \max_{1 \leq i \leq (1-\alpha)t} \eta_i^{(t)} \max_{1 \leq i \leq (1-\alpha)t} \left(\|\bar{P}^{\pi_{t-1}} \Delta_{i-1}\|_\infty + \|\Lambda \hat{P}_{\mathcal{K}}^{(t)} \mathbf{V}_{i-1}\|_\infty + \|\Lambda P_{\mathcal{K}} \mathbf{V}_{i-1}\|_\infty \right) \\
&\stackrel{(i)}{\leq} \eta_0^{(t)} \|\Delta_0\|_\infty + t \max_{1 \leq i \leq (1-\alpha)t} \eta_i^{(t)} \max_{1 \leq i \leq (1-\alpha)t} (\|\Delta_{i-1}\|_\infty + 2 \|\mathbf{V}_{i-1}\|_\infty) \\
&\stackrel{(ii)}{\leq} \frac{1}{2T^2} \cdot \frac{1}{1-\gamma} + \frac{1}{2T^2} \cdot t \cdot \frac{3}{1-\gamma} \\
&\leq \frac{2}{(1-\gamma)T},
\end{aligned} \tag{46}$$

where (i) is due to the fact that $\|\bar{P}^{\pi_{t-1}}\|_1 = \|\Lambda \hat{P}_{\mathcal{K}}^{(t)}\|_1 = \|\Lambda P_{\mathcal{K}}\|_1 = 1$ and (ii) comes from [4, Equation (39a)].

Bounding ν_t . To control the second term, we apply the following Freedman's inequality.

Lemma 3 (Freedman's Inequality). *Consider a real-valued martingale $\{Y_k : k = 0, 1, 2, \dots\}$ with difference sequence $\{X_k : k = 1, 2, 3, \dots\}$. Assume that the difference sequence is uniformly bounded:*

$$|X_k| \leq R \quad \text{and} \quad \mathbb{E}[X_k | \{X_j\}_{j=1}^{k-1}] = 0 \quad \text{for all } k \geq 1.$$

Let

$$S_n := \sum_{k=1}^n X_k, \quad T_n := \sum_{k=1}^n \text{Var}\{X_k | \{X_j\}_{j=1}^{k-1}\}.$$

Then for any given $\sigma^2 \geq 0$, one has

$$\mathbb{P}(|S_n| \geq \tau \text{ and } T_n \leq \sigma^2) \leq 2 \exp\left(-\frac{\tau^2/2}{\sigma^2 + R\tau/3}\right).$$

In addition, suppose that $W_n \leq \sigma^2$ holds deterministically. For any positive integer $K \geq 1$, with probability at least $1 - \delta$ one has

$$|S_n| \leq \sqrt{8 \max\left\{T_n, \frac{\sigma^2}{2K}\right\} \log \frac{2K}{\delta}} + \frac{4}{3} R \log \frac{2K}{\delta}.$$

Proof. See [4, Theorem 4]. □

To apply this inequality, we can express ν_t as

$$\nu_t := \sum_{i=(1-\alpha)t+1}^t \mathbf{x}_i,$$

with

$$\mathbf{x}_i := \eta_i^{(t)} \gamma \Lambda \left(\hat{P}_{\mathcal{K}}^{(t)} - P_{\mathcal{K}} \right) \mathbf{V}_{i-1}, \quad \text{and} \quad \mathbb{E}[\mathbf{x}_i | \mathbf{V}_{i-1}, \dots, \mathbf{V}_0] = \mathbf{0}. \tag{47}$$

1. In order to calculate bound R in Lemma 3, one has

$$\begin{aligned}
B &:= \max_{(1-\alpha)t < t \leq t} \|\mathbf{x}_i\|_\infty \leq \max_{(1-\alpha)t < t \leq t} \left\| \eta_i^{(t)} \Lambda \left(\hat{P}_{\mathcal{K}}^{(t)} - P_{\mathcal{K}} \right) \mathbf{V}_{i-1} \right\|_\infty \\
&\leq \max_{(1-\alpha)t < t \leq t} \eta_i^{(t)} \left(\|\Lambda \hat{P}_{\mathcal{K}}^{(t)}\|_1 + \|\Lambda P_{\mathcal{K}}\|_1 \right) \|\mathbf{V}_{i-1}\|_\infty \\
&\leq \max_{(1-\alpha)t < t \leq t} \eta_i^{(t)} \cdot \frac{2}{1-\gamma} \leq \frac{4 \log^4 T}{(1-\gamma)^2 T},
\end{aligned}$$

where the last inequality comes from [4, Eqn (39b)] and the fact that $\|\mathbf{V}_{i-1}\|_\infty \leq \frac{1}{1-\gamma}$.

2. Then regarding the variance term, we claim for the moment that

$$\begin{aligned}\mathbf{W}_t &:= \sum_{i=(1-\alpha)t+1}^t \text{diag}(\text{Var}(\mathbf{x}_i | \mathbf{V}_{i-1}, \dots, \mathbf{V}_0)) \\ &\leq \gamma^2 \sum_{i=(1-\alpha)t+1}^t \left(\eta_i^{(t)}\right)^2 \text{Var}_{\overline{\mathbf{P}}}(\mathbf{V}_{i-1}).\end{aligned}\quad (48)$$

Then we have

$$\begin{aligned}\mathbf{W}_t &\leq \max_{(1-\alpha)t \leq i \leq t} \eta_i^{(t)} \left(\sum_{i=(1-\alpha)t+1}^t \eta_i^{(t)} \right) \max_{(1-\alpha)t \leq i < t} \text{Var}_{\overline{\mathbf{P}}}(\mathbf{V}_i) \\ &\leq \frac{2 \log^4 T}{(1-\gamma)T} \max_{(1-\alpha)t \leq i < t} \text{Var}_{\overline{\mathbf{P}}}(\mathbf{V}_i),\end{aligned}\quad (49)$$

where the second line comes from [4, Eqns (39b), (40)]. A trivial upper bound for \mathbf{W}_t is

$$|\mathbf{W}_t| \leq \frac{2 \log^4 T}{(1-\gamma)T} \cdot \frac{1}{(1-\gamma)^2} \mathbf{1} = \frac{2 \log^4 T}{(1-\gamma)^3 T} \mathbf{1},$$

which uses the fact that $\text{Var}_{\mathbf{P}}(\mathbf{V}_i) \leq \|\mathbf{V}_i\|_\infty^2 \leq 1/(1-\gamma)^2$.

Then, we invoke Lemma 3 with $K = \left\lceil 2 \log_2 \frac{1}{1-\gamma} \right\rceil$ and apply the union bound argument over \mathcal{K} to arrive at

$$\begin{aligned}|\boldsymbol{\nu}_t| &\leq \sqrt{8 \left(\mathbf{W}_t + \frac{\sigma^2}{2^K} \mathbf{1} \right) \log \frac{8KT \log \frac{1}{1-\gamma}}{\delta}} + \frac{4}{3} B \log \frac{8KT \log \frac{1}{1-\gamma}}{\delta} \mathbf{1} \\ &\leq \sqrt{8 \left(\mathbf{W}_t + \frac{2 \log^4 T}{(1-\gamma)T} \mathbf{1} \right) \log \frac{8KT}{\delta}} + \frac{4}{3} B \log \frac{8KT \log \frac{1}{1-\gamma}}{\delta} \mathbf{1} \\ &\leq \sqrt{\frac{32 \log^4 T}{(1-\gamma)T} \log \frac{8KT}{\delta} \left(\max_{(1-\alpha)t \leq i < t} \text{Var}_{\Lambda \mathbf{P}_{\mathcal{K}}}(\mathbf{V}_i) + \mathbf{1} \right)} + \frac{12 \log^4 T}{(1-\gamma)^2 T} \log \frac{8KT}{\delta} \mathbf{1}.\end{aligned}\quad (50)$$

Hence if we define

$$\boldsymbol{\varphi}_t := 64 \frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)T} \left(\max_{\frac{t}{2} \leq i \leq t} \text{Var}_{\overline{\mathbf{P}}}(\mathbf{V}_i) + \mathbf{1} \right),$$

then (46) and (50) implies that

$$|\boldsymbol{\theta}_t| + |\boldsymbol{\nu}_t| + |\boldsymbol{\omega}_t| \leq \sqrt{\boldsymbol{\varphi}_t} + \frac{2\gamma\xi}{1-\gamma} \mathbf{1},\quad (51)$$

with probability over $1-\delta$ for all $2t/3 \leq k \leq t$, as long as $T \gg \log^4 T \log \frac{KT}{\delta} / (1-\gamma)^3$. Therefore, plugging (51) into (45), we arrive at the recursive relationship

$$\boldsymbol{\Delta}_t \leq \sqrt{\boldsymbol{\varphi}_t} + \frac{2\gamma\xi}{1-\gamma} \mathbf{1} + \sum_{i=(1-\alpha)k+1}^k \eta_i^{(k)} \gamma \overline{\mathbf{P}}^{\pi_{i-1}} \boldsymbol{\Delta}_{i-1} = \sqrt{\boldsymbol{\varphi}_t} + \frac{2\gamma\xi}{1-\gamma} \mathbf{1} + \sum_{i=(1-\alpha)k}^{k-1} \eta_i^{(k)} \gamma \overline{\mathbf{P}}^{\pi_{i-1}} \boldsymbol{\Delta}_i.$$

This recursion is expressed in a similar way as [4, Eqn. (46)] so we can invoke similar derivation in [4, Appendix A.2] to obtain that

$$\boldsymbol{\Delta}_t \leq 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} \left(1 + \max_{\frac{t}{2} \leq i < t} \|\boldsymbol{\Delta}_i\|_\infty \right)} \mathbf{1} + \frac{2\gamma\xi}{(1-\gamma)^2} \mathbf{1}.\quad (52)$$

Then we turn to (44). Applying a similar argument, we can deduce that

$$\boldsymbol{\Delta}_t \geq -30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} \left(1 + \max_{\frac{t}{2} \leq i < t} \|\boldsymbol{\Delta}_i\|_\infty \right)} \mathbf{1} - \frac{2\gamma\xi}{(1-\gamma)^2} \mathbf{1}.\quad (53)$$

For any t satisfying $\frac{T}{c_2 \log \frac{1}{1-\gamma}} \leq t \leq T$, taking (52) and (53) collectively gives rise to

$$\|\Delta_t\|_\infty \leq 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} \left(1 + \max_{\frac{t}{2} \leq i < t} \|\Delta_i\|_\infty\right)} + \frac{2\gamma\xi}{(1-\gamma)^2}. \quad (54)$$

Let

$$u_k := \max \left\{ \|\Delta_t\|_\infty : 2^k \frac{T}{c_2 \log \frac{1}{1-\gamma}} \leq t \leq T \right\}.$$

By taking supremum over $t \in \{[2^k T / (c_2 \log \frac{1}{1-\gamma})], \dots, T\}$ on both sides of (54), we have

$$u_k \leq 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} (1 + u_{k-1})} + \frac{2\gamma\xi}{(1-\gamma)^2} \quad \forall 1 \leq k \leq \log \left(c_2 \log \frac{1}{1-\gamma} \right). \quad (55)$$

It is straightforward to bound $u_0 \leq \frac{1}{1-\gamma}$. For $k \geq 1$, it is straightforward to obtain from (55) that

$$u_k \leq 3 \max \left\{ 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}}, 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} u_{k-1}}, \frac{2\gamma\xi}{(1-\gamma)^2} \right\}, \quad (56)$$

for $1 \leq k \leq \log(c_2 \log \frac{1}{1-\gamma})$. We analyze (56) under two different cases:

1. If there exists some integer k_0 with $1 \leq k_0 < \lceil \log(c_2 \log \frac{1}{1-\gamma}) \rceil$, such that

$$u_{k_0} \leq \max \left\{ 1, \frac{6\gamma\xi}{(1-\gamma)^2} \right\},$$

then it is straightforward to check from (56) that

$$u_{k_0+1} \leq 3 \max \left\{ 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}}, \frac{2\gamma\xi}{(1-\gamma)^2} \right\} \quad (57)$$

as long as $T \geq C_3(1-\gamma)^{-4} \log^4 T \log(KT/\delta)$ for some sufficiently large constant $C_3 > 0$.

2. Otherwise we have $u_k > \max\{1, \frac{6\gamma\xi}{(1-\gamma)^2}\}$ for all $1 \leq k < \lceil \log(c_2 \log \frac{1}{1-\gamma}) \rceil$. This together with (56) suggests that

$$\max \left\{ 1, \frac{6\gamma\xi}{(1-\gamma)^2} \right\} < 3 \max \left\{ 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}}, 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} u_{k-1}}, \frac{2\gamma\xi}{(1-\gamma)^2} \right\},$$

and therefore

$$\max \left\{ 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}}, 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} u_{k-1}}, \frac{2\gamma\xi}{(1-\gamma)^2} \right\} = 30 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} u_{k-1}}$$

for all $1 \leq k \leq \log(c_2 \log \frac{1}{1-\gamma})$. Let

$$v_k := 90 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} u_{k-1}}.$$

Then we know from (55) that

$$u_k \leq v_k \quad \forall 1 \leq k \leq \log \left(c_2 \log \frac{1}{1-\gamma} \right).$$

By applying the above two inequalities recursively, we know that

$$\begin{aligned}
u_k &\leq v_k = \left(\frac{8100 \log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} \right)^{1/2} u_{k-1}^{1/2} \leq \left(\frac{8100 \log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} \right)^{1/2} v_{k-1}^{1/2} \\
&\leq \left(\frac{8100 \log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} \right)^{1/2+1/4} u_{k-2}^{1/4} \leq \left(\frac{8100 \log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} \right)^{1/2+1/4} v_{k-2}^{1/4} \\
&\leq \dots \leq \left(\frac{8100 \log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T} \right)^{1-1/2^k} u_0^{1/2^k} \leq \sqrt{\frac{8100 \log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}} \left(\frac{1}{1-\gamma} \right)^{1/2^k},
\end{aligned}$$

where the last inequality holds as long as $T \geq C_3 \log^4 T \log(KT/\delta)(1-\gamma)^{-4}$ for some sufficiently large constant $C_3 > 0$. Let $k_0 = \tilde{c} \log \log \frac{1}{1-\gamma}$ for some properly chosen constant $\tilde{c} > 0$ such that k_0 is an integer between 1 and $\log(c_2 \log \frac{1}{1-\gamma})$, we have

$$u_{k_0} \leq \sqrt{\frac{8100 \log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}} \left(\frac{1}{1-\gamma} \right)^{1/2^{k_0}} = O \left(\sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}} \right).$$

When $T \geq C_3 \log^4 T \log(KT/\delta)(1-\gamma)^{-4}$ for some sufficiently large constant $C_3 > 0$, this implies that $u_{k_0} < 1$, which contradicts with the preassumption that $u_k > \max\{1, \frac{6\gamma\xi}{(1-\gamma)^2}\}$ for all $1 \leq k \leq c_2 \log \frac{1}{1-\gamma}$.

Consequently, (57) must hold true and then the definition of u_k immediately leads to

$$\|\Delta_T\|_\infty \leq 90 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}} + \frac{6\gamma\xi}{(1-\gamma)^2}.$$

Then for any $\varepsilon \in (0, 1]$, one has

$$\|\Delta_T\|_\infty \leq \varepsilon + \frac{6\gamma\xi}{(1-\gamma)^2},$$

as long as

$$90 \sqrt{\frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 T}} \leq \varepsilon.$$

Hence, if the total number of iterations T satisfies

$$T \geq C_3 \frac{\log^4 T \log \frac{KT}{\delta}}{(1-\gamma)^4 \varepsilon^2}$$

for some sufficiently large constant $C_3 > 0$, (10) would hold for Algorithm 1 with probability over $1 - \delta$.

Finally, we are left to justify (48). Recall the definition of \mathbf{x}_i (cf. (47)), one has

$$\begin{aligned}
\text{diag}(\text{Var}(\mathbf{x}_i | \mathbf{V}_{i-1}, \dots, \mathbf{V}_0)) &= \gamma^2 \left(\eta_i^{(t)} \right)^2 \text{diag} \left(\text{Var} \left(\Lambda \left(\hat{\mathbf{P}}_{\mathcal{K}}^{(t)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{i-1} | \mathbf{V}_{i-1} \right) \right) \\
&= \gamma^2 \left(\eta_i^{(t)} \right)^2 \text{diag} \left(\Lambda \text{Var} \left(\left(\hat{\mathbf{P}}_{\mathcal{K}}^{(i)} - \mathbf{P}_{\mathcal{K}} \right) \mathbf{V}_{i-1} | \mathbf{V}_{i-1} \right) \Lambda^\top \right) \\
&= \gamma^2 \left(\eta_i^{(t)} \right)^2 \left\{ \lambda(s, a)^2 \text{Var}_{\mathbf{P}_{\mathcal{K}}}(\mathbf{V}_{i-1}) \right\}_{s,a},
\end{aligned}$$

where the notation $\text{Var}_{\mathbf{P}_{\mathcal{K}}}(\mathbf{V}_{i-1})$ is defined in (12). Plugging this into the definition of \mathbf{W}_t leads to

$$\begin{aligned}
\mathbf{W}_t &= \gamma^2 \sum_{i=(1-\alpha)t+1}^t \left(\eta_i^{(t)} \right)^2 \left\{ \lambda(s, a)^2 \text{Var}_{\mathbf{P}_{\mathcal{K}}}(\mathbf{V}_{i-1}) \right\}_{s,a} \\
&= \gamma^2 \sum_{i=(1-\alpha)t+1}^t \left(\eta_i^{(t)} \right)^2 \left\{ \lambda(s, a)^2 (\mathbf{P}_{\mathcal{K}}(\mathbf{V}_{i-1} \circ \mathbf{V}_{i-1}) - (\mathbf{P}_{\mathcal{K}} \mathbf{V}_{i-1}) \circ (\mathbf{P}_{\mathcal{K}} \mathbf{V}_{i-1})) \right\}_{s,a}. \quad (58)
\end{aligned}$$

Then we introduce a useful claim as follows. The proof is deferred to Appendix C.2.

Claim 1. For any state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$ and vector $\mathbf{V} \in \mathbb{R}^{|\mathcal{S}|}$, one has

$$\begin{aligned} & \lambda(s, a)^2 (\mathbf{P}_{\mathcal{K}}(\mathbf{V} \circ \mathbf{V}) - (\mathbf{P}_{\mathcal{K}}\mathbf{V}) \circ (\mathbf{P}_{\mathcal{K}}\mathbf{V})) \\ & \leq \lambda(s, a) \mathbf{P}_{\mathcal{K}}(\mathbf{V} \circ \mathbf{V}) - (\lambda(s, a) \mathbf{P}_{\mathcal{K}}\mathbf{V}) \circ (\lambda(s, a) \mathbf{P}_{\mathcal{K}}\mathbf{V}). \end{aligned} \quad (59)$$

By invoking this claim with $\mathbf{V} = \mathbf{V}^{i-1}$ and taking collectively with (58), one has

$$\begin{aligned} \mathbf{W}_t & \leq \gamma^2 \sum_{i=(1-\beta)t+1}^t \left(\eta_i^{(t)} \right)^2 \{ \lambda(s, a) \mathbf{P}_{\mathcal{K}}(\mathbf{V}_{i-1} \circ \mathbf{V}_{i-1}) - (\lambda(s, a) \mathbf{P}_{\mathcal{K}}\mathbf{V}_{i-1}) \circ (\lambda(s, a) \mathbf{P}_{\mathcal{K}}\mathbf{V}_{i-1}) \}_{s,a} \\ & = \gamma^2 \sum_{i=(1-\beta)t+1}^t \left(\eta_i^{(t)} \right)^2 [\Lambda \mathbf{P}_{\mathcal{K}}(\mathbf{V}_{i-1} \circ \mathbf{V}_{i-1}) - (\Lambda \mathbf{P}_{\mathcal{K}}\mathbf{V}_{i-1}) \circ (\Lambda \mathbf{P}_{\mathcal{K}}\mathbf{V}_{i-1})] \\ & = \gamma^2 \sum_{i=(1-\beta)t+1}^t \left(\eta_i^{(t)} \right)^2 \text{Var}_{\overline{\mathbf{P}}}(\mathbf{V}_{i-1}), \end{aligned}$$

which is the desired result.

C.2 Proof of Claim 1

To simplify notations in this proof, we use $[\lambda_i]_{i=1}^K$, $[P_{i,j}]_{1 \leq i \leq K, 1 \leq j \leq |\mathcal{S}|}$ and $[V_i]_{i=1}^{|\mathcal{S}|}$ to denote $\lambda(s, a)$, $\mathbf{P}_{\mathcal{K}}$ and \mathbf{V} respectively. Then one has

$$\begin{aligned} & \lambda(s, a) \mathbf{P}_{\mathcal{K}}(\mathbf{V} \circ \mathbf{V}) - (\lambda(s, a) \mathbf{P}_{\mathcal{K}}\mathbf{V}) \circ (\lambda(s, a) \mathbf{P}_{\mathcal{K}}\mathbf{V}) \\ & - \lambda(s, a)^2 (\mathbf{P}_{\mathcal{K}}(\mathbf{V} \circ \mathbf{V}) - (\mathbf{P}_{\mathcal{K}}\mathbf{V}) \circ (\mathbf{P}_{\mathcal{K}}\mathbf{V})) \\ & = \sum_{i=1}^K \sum_{j=1}^{|\mathcal{S}|} \lambda_i P_{i,j} V_j^2 - \left(\sum_{i=1}^K \sum_{j=1}^{|\mathcal{S}|} \lambda_i P_{i,j} V_j \right)^2 - \sum_{i=1}^K \sum_{j=1}^{|\mathcal{S}|} \lambda_i^2 P_{i,j} V_j^2 + \sum_{i=1}^K \lambda_i^2 \left(\sum_{j=1}^{|\mathcal{S}|} P_{i,j} V_j \right)^2 \\ & = \sum_{i=1}^K \sum_{j=1}^{|\mathcal{S}|} \lambda_i P_{i,j} V_j \left[(1 - \lambda_i) V_j - \sum_{i' \neq i} \sum_{j'=1}^{|\mathcal{S}|} \lambda_{i'} P_{i',j'} V_{j'} \right] \\ & = \sum_{i=1}^K \sum_{j=1}^{|\mathcal{S}|} \lambda_i P_{i,j} V_j \left[\left(\sum_{i'=1}^K \sum_{j'=1}^{|\mathcal{S}|} \lambda_{i'} P_{i',j'} - \lambda_i \right) V_j - \sum_{i' \neq i} \sum_{j'=1}^{|\mathcal{S}|} \lambda_{i'} P_{i',j'} V_{j'} \right] \\ & = \sum_{i=1}^K \sum_{j=1}^{|\mathcal{S}|} \sum_{i' \neq i} \sum_{j'=1}^{|\mathcal{S}|} \lambda_i P_{i,j} V_j \lambda_{i'} P_{i',j'} (V_j - V_{j'}) \end{aligned}$$

where in the penultimate equality, we use the fact that

$$\sum_{i'=1}^K \sum_{j'=1}^{|\mathcal{S}|} \lambda_{i'} P_{i',j'} = \lambda(s, a) \mathbf{P}_{\mathcal{K}} \mathbf{1} = 1.$$

It follows that

$$\begin{aligned}
& \lambda(s, a) \mathbf{P}_{\mathcal{K}}(\mathbf{V} \circ \mathbf{V}) - (\lambda(s, a) \mathbf{P}_{\mathcal{K}} \mathbf{V}) \circ (\lambda(s, a) \mathbf{P}_{\mathcal{K}} \mathbf{V}) \\
& \quad - \lambda(s, a)^2 (\mathbf{P}_{\mathcal{K}}(\mathbf{V} \circ \mathbf{V}) - (\mathbf{P}_{\mathcal{K}} \mathbf{V}) \circ (\mathbf{P}_{\mathcal{K}} \mathbf{V})) \\
& = \sum_{i=1}^K \sum_{1 \leq i' < i} \sum_{j=1}^{|\mathcal{S}|} \sum_{j'=1}^{|\mathcal{S}|} [\lambda_i P_{i,j} V_j \lambda_{i'} P_{i',j'} (V_j - V_{j'}) + \lambda_{i'} P_{i',j} V_j \lambda_i P_{i,j'} (V_j - V_{j'})] \\
& = \sum_{i=1}^K \sum_{1 \leq i' < i} \lambda_i \lambda_{i'} \left[\sum_{j=1}^{|\mathcal{S}|} \sum_{j'=1}^{|\mathcal{S}|} P_{i,j} V_j P_{i',j'} (V_j - V_{j'}) + \sum_{j=1}^{|\mathcal{S}|} \sum_{j'=1}^{|\mathcal{S}|} P_{i',j} V_j P_{i,j'} (V_j - V_{j'}) \right] \\
& \stackrel{(i)}{=} \sum_{i=1}^K \sum_{1 \leq i' < i} \lambda_i \lambda_{i'} \left[\sum_{j=1}^{|\mathcal{S}|} \sum_{j'=1}^{|\mathcal{S}|} P_{i,j} V_j P_{i',j'} (V_j - V_{j'}) + \sum_{j=1}^{|\mathcal{S}|} \sum_{j'=1}^{|\mathcal{S}|} P_{i',j'} V_{j'} P_{i,j} (V_{j'} - V_j) \right] \\
& = \sum_{i=1}^K \sum_{1 \leq i' < i} \lambda_i \lambda_{i'} \left[\sum_{j=1}^{|\mathcal{S}|} \sum_{j'=1}^{|\mathcal{S}|} P_{i,j} P_{i',j'} (V_j - V_{j'})^2 \right] \\
& \geq 0,
\end{aligned}$$

where in (i), we exchange the indices j and j' .

D Feature dimension and the number of anchor state-action pairs

The assumption that the feature dimension (denoted by K_d) and the number of anchor state-action pairs (denoted by K_n) are equal is actually non-essential. In what follows, we will show that if $K_d \neq K_n$, then we can modify the current feature mapping $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^{K_d}$ to achieve a new feature mapping $\phi' : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^{K_n}$ that does not change the transition model P . By doing so, the new feature dimension K_n equals to the number of anchor state-action pairs.

To begin with, we recall from Definition 1 that there exists K_d unknown functions $\psi_1, \dots, \psi_{K_d} : \mathcal{S} \rightarrow \mathbb{R}$, such that

$$P(s'|s, a) = \sum_{k=1}^{K_d} \phi_k(s, a) \psi_k(s'),$$

for every $(s, a) \in \mathcal{S} \times \mathcal{A}$ and $s' \in \mathcal{S}$. In addition, we also recall from Assumption 1 that there exists $\mathcal{K} \subseteq \mathcal{S} \times \mathcal{A}$ with $|\mathcal{K}| = K_n$ such that for any $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\phi(s, a) = \sum_{i: (s_i, a_i) \in \mathcal{K}} \lambda_i(s, a) \phi(s_i, a_i) \in \mathbb{R}^{K_d} \quad \text{for} \quad \sum_{i=1}^{K_n} \lambda_i(s, a) = 1 \quad \text{and} \quad \lambda_i(s, a) \geq 0.$$

Case 1: $K_d > K_n$. In this case, the vectors in $\{\phi(s, a) : (s, a) \in \mathcal{K}\}$ are linearly independent. For ease of presentation and without loss of generality, we assume that $K_d = K_n + 1$. This indicates that the matrix $\Phi \in \mathbb{R}^{K_d \times (|\mathcal{S}| |\mathcal{A}|)}$ whose columns are composed of the feature vectors of all state-action pairs has rank K_n and is hence not full row rank. This suggests that there exists K_n linearly independent rows (without loss of generality, we assume they are the first K_n rows). We can remove the last row from Φ to obtain $\Phi' := \Phi_{1:K_n, :} \in \mathbb{R}^{K_n \times (|\mathcal{S}| |\mathcal{A}|)}$ such that Φ' is full row rank. Then we show that we can actually use the columns of Φ' as new feature mappings. To see why this is true, note that the last row $\Phi_{K_n+1, :}$ can be represented as a linear combination of the first K_n rows, namely there must exist constants $\{c_k\}_{k=1}^{K_n}$ such that for any $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\phi_{K_n+1}(s, a) = \sum_{k=1}^{K_n} c_k \phi_k(s, a).$$

Define $\psi'_k = \psi_k + c_k \psi_{K_n+1}$ for $k = 1, \dots, K_n$, we have

$$\begin{aligned} P(s'|s, a) &= \sum_{k=1}^{K_d} \phi_k(s, a) \psi_k(s') = \phi_{K_n+1}(s, a) \psi_{K_n+1}(s') + \sum_{k=1}^{K_n} \phi_k(s, a) \psi_k(s') \\ &= \sum_{k=1}^{K_n} \phi_k(s, a) [\psi_k(s') + c_k \psi_{K_n+1}(s')] = \sum_{k=1}^{K_n} \phi_k(s, a) \psi'_k(s'), \end{aligned}$$

which is linear with respect to the new K_n dimensional feature vectors. It is also straightforward to check that the new feature mapping satisfies Assumption 1 with the original anchor state-action pairs \mathcal{K} .

Case 2: $K_d < K_n$. For ease of presentation and without loss of generality, we assume that $K_n = K_d + 1$ and that the subspace spanned by the feature vectors of anchor state-action pairs is non-degenerate, i.e., has rank K_d (otherwise we can use similar method as in Case 1 to further reduce the feature dimension K_d). In this case, the matrix $\Phi_{\mathcal{K}} \in \mathbb{R}^{K_d \times K_n}$ whose columns are composed of the feature vectors of anchor state-action pairs has rank K_d . We can add $K_n - K_d = 1$ new row to $\Phi_{\mathcal{K}}$ to obtain $\Phi'_{\mathcal{K}} \in \mathbb{R}^{K_n \times K_n}$ such that $\Phi'_{\mathcal{K}}$ has full rank K_n . Then we let the columns of $\Phi'_{\mathcal{K}} = [\phi'(s, a)]_{(s,a) \in \mathcal{K}}$ to be the new feature vectors of the anchor state-action pairs, and define the new feature vectors for all other state-action pairs $(s, a) \notin \mathcal{K}$ by

$$\phi'(s, a) = \sum_{i: (s_i, a_i) \in \mathcal{K}} \lambda_i(s, a) \phi'(s_i, a_i).$$

We can check that the transition model P is not changed if we let $\psi_{K_n}(s') = 0$ for every $s' \in \mathcal{S}$. It is also straightforward to check that Assumption 1 is satisfied.

To conclude, when $K_d \neq K_n$, we can always construct a new set of feature mappings with dimension K_n such that: (i) the feature dimension equals to the number of anchor state-action pairs (they are both K_n); (ii) the transition model can still be linearly parameterized by this new set of feature mappings; and (iii) the anchor state-action pair assumption (Assumption 1) is satisfied with the original anchor state-action pairs.

References

- [1] Alekh Agarwal, Sham Kakade, and Lin F Yang. Model-based reinforcement learning with a generative model is minimax optimal. In *Conference on Learning Theory*, pages 67–83. PMLR, 2020.
- [2] Sanjeev Arora, Rong Ge, and Ankur Moitra. Learning topic models—going beyond svd. In *2012 IEEE 53rd annual symposium on foundations of computer science*, pages 1–10. IEEE, 2012.
- [3] Mohammad Gheshlaghi Azar, Rémi Munos, and Hilbert J Kappen. Minimax pac bounds on the sample complexity of reinforcement learning with a generative model. *Machine learning*, 91(3):325–349, 2013.
- [4] Gen Li, Changxiao Cai, Yuxin Chen, Yuantao Gu, Yuting Wei, and Yuejie Chi. Is q-learning minimax optimal? a tight sample complexity analysis. *arXiv preprint arXiv:2102.06548*, 2021.
- [5] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge University Press, 2018.